

Catalytic Prior Distributions with Application to Generalized Linear Models

Dongming Huang^a, Nathan Stein^b, Donald B. Rubin^{a,c,1}, and S. C. Kou^{a,1}

^aDepartment of Statistics, Harvard University, Cambridge, MA 02138; ^bSpotify, New York, NY 10011; ^cTsinghua University, Beijing 100084, China

This manuscript was compiled on April 2, 2020

1 A catalytic prior distribution is designed to stabilize a high-
2 dimensional “working model” by shrinking it toward a “simplified
3 model.” The shrinkage is achieved by supplementing the observed
4 data with a small amount of “synthetic data” generated from a predic-
5 tive distribution under the simpler model. We apply this framework
6 to generalized linear models, where we propose various strategies
7 for the specification of a tuning parameter governing the degree of
8 shrinkage and study resultant theoretical properties. In simulations,
9 the resulting posterior estimation using such a catalytic prior outper-
10 forms maximum likelihood estimation from the working model and
11 is generally comparable or superior to existing competitive methods
12 in terms of frequentist prediction accuracy of point estimation and
13 coverage accuracy of interval estimation. The catalytic priors have
14 simple interpretations and are easy to formulate.

Bayesian priors | synthetic-data | stable estimation | predictive distribution | regularization

1 The prior distribution is a unique and important feature
2 of Bayesian analysis, yet in practice, it can be difficult
3 to quantify existing knowledge into actual prior distributions;
4 thus, automated construction of prior distributions can be
5 desirable. Such prior distributions should stabilize posterior
6 estimation in situations when maximum likelihood (ML) be-
7 haves problematically, which can occur when sample sizes are
8 small relative to the dimensionality of the models. Here we
9 propose a class of prior distributions designed to address such
10 situations. Henceforth we call the complex model that the
11 investigator wishes to use to analyze the data the “working
12 model.”

13 Often with real working models and datasets, the sample
14 sizes are relatively small, and a likelihood-based analysis is un-
15 stable, whereas a likelihood-based analysis of the same dataset
16 using a simpler but less rich model can be stable. Catalytic
17 priors* effectively supplement the observed data with a small
18 amount of *synthetic data* generated from a suitable predictive
19 distribution, such as the posterior predictive distribution under
20 the simpler model. In this way, the resulting posterior distri-
21 bution under the working model is pulled toward the posterior
22 distribution under the simpler model, resulting in estimates
23 and predictions with better frequentist properties. The name
24 for these priors arises because a catalyst is something that
25 stimulates a reaction to take place that would not take place
26 (or not as effectively) without it, but only an insubstantial
27 amount of the catalyst is needed. When the information in the
28 observed data is substantial, the catalytic prior has a minor
29 influence on the resulting inference, because the information
30 in the synthetic data is small relative to the information in
31 the observed data.

32 We are not the first to suggest such priors, but we embed

33 the suggestion within a general framework designed for a broad
34 range of examples. One early suggestion for the applied use
35 of such priors was in Ref. (1), which was based on an earlier
36 proposal by Rubin in a 1983 report for the U.S. Census Bureau
37 (reprinted as an appendix in Ref. (2)). Such a prior was also
38 used in a Bayesian analysis of data with noncompliance in a
39 randomized trial (3).

As in both of these earlier references, consider logistic
regression as an example:

$$y_i | \mathbf{x}_i, \beta \sim \text{Bernoulli}\left(1/(1 + \exp(-\mathbf{x}_i^\top \beta))\right), \quad i = 1, \dots, n,$$

40 where, for the i th data point (y_i, \mathbf{x}_i) , $y_i \in \{0, 1\}$ is the re-
41 sponse, and $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip-1})^\top$ represents p covariates,
42 with unknown coefficients $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})^\top$. The ML
43 estimate (MLE) of β is infinite when there is complete separa-
44 tion (4, 5) of the observed covariate values in the two
45 response categories, which can occur easily when p is large
46 relative to n . Earlier attempts to address this problem, such
47 as using Jeffrey’s prior (6–9), are not fully satisfactory. This
48 problem arises commonly in practice, for example, Ref. (1)
49 studied the mapping of industry and occupation (I/O) codes
50 in the 1970 U.S. Census to the 1980 Census codes, where
51 both coding systems had hundreds of categories. The I/O
52 classification system changed drastically from the 1970 Census
53 to the 1980 Census, and a single 1970 code could map into
54 as many as 60 possible 1980 codes. For each 1970 code, the
55 1980 code was considered as missing and multiply-imputed
56 based on covariates. The imputation models were nested (di-
57 chotomous) logistic regression models (10) estimated from a
58 special training sample for which both 1970 and 1980 codes
59 were known. The covariates used in these models were derived
60 from nine different factors (sex, age, race, etc.) that formed
61 a cross-classification with $J = 2,304$ categories. The sample
62 available to estimate the mapping was smaller than ten for
63 some 1970 codes, and many of these logistic regression models

Significance Statement

We propose a strategy for building prior distributions that stabilize the estimation of complex “working models” when sample sizes are too small for standard statistical analysis. The stabilization is achieved by supplementing the observed data with a small amount of synthetic data generated from the predictive distribution of a simpler model. This class of prior distributions is easy to use and allows direct statistical interpretation.

D.H., N.S., D.B.R. and S.C.K. designed research, performed research, analyzed data, and wrote the paper.

The authors declare no conflict of interest.

*Throughout the paper, we use the inelegant but compact ‘priors’ in place of the correct ‘prior distributions’

¹To whom correspondence should be addressed. E-mail: kou@stat.harvard.edu or dbrubin@me.com

64 faced complete separation. The successful approach in Ref.
65 (1) was to use the prior distribution

$$66 \quad \pi(\beta) \propto \prod_{j=1}^J \left(\frac{e^{\mathbf{x}_j^* \top \beta}}{1 + e^{\mathbf{x}_j^* \top \beta}} \right)^{p\hat{\mu}/J} \left(\frac{1}{1 + e^{\mathbf{x}_j^* \top \beta}} \right)^{p(1-\hat{\mu})/J}, \quad [1]$$

where each \mathbf{x}_j^* is a possible covariate vector of the cross-classification; p is the dimension of β ; and $\hat{\mu} = \sum_{i=1}^n y_i/n$ is the marginal proportion of ones among the observed responses. In this example, the simpler model has the responses y_i 's independent of the covariates:

$$y_i \mid \mathbf{x}_i, \mu \sim \text{Bernoulli}(\mu) \quad (i = 1, \dots, n),$$

where $\mu \in (0, 1)$ is a probability estimated by $\hat{\mu}$. If we supplement the dataset with $p\hat{\mu}/J$ synthetic data points ($y_j^* = 1, \mathbf{x}_j^*$) and $p(1 - \hat{\mu})/J$ synthetic data points ($y_j^* = 0, \mathbf{x}_j^*$) for each \mathbf{x}_j^* ($j = 1, \dots, J$), then the likelihood function of the augmented dataset has the same form as the posterior distribution with the prior in Eq. (1):

$$\pi(\beta \mid \{(y_i, \mathbf{x}_i)\}_{i=1}^n) \quad [2] \\ \propto \prod_{j=1}^J \left(\frac{e^{\mathbf{x}_j^* \top \beta}}{1 + e^{\mathbf{x}_j^* \top \beta}} \right)^{N_{j,1} + p\hat{\mu}/J} \left(\frac{1}{1 + e^{\mathbf{x}_j^* \top \beta}} \right)^{N_{j,0} + p(1-\hat{\mu})/J},$$

67 where $N_{j,1}, N_{j,0}$ are, respectively, the numbers of $(1, \mathbf{x}_j^*)$ and
68 $(0, \mathbf{x}_j^*)$ in the observed data. In this construction, the total
69 amount of synthetic data is taken to be p , the dimension of
70 β ; see Remark 2.2 in SI for more discussion. The resulting
71 MLE with the augmented dataset equals the maximum posterior
72 estimator (the value of β that maximizes the posterior
73 distribution), and it will always be unique and finite when
74 $\hat{\mu} \in (0, 1)$.

75 How to use the synthetic-data perspective for constructing
76 general prior distributions, which we called catalytic prior
77 distributions, is our focus. We mathematically formulate the
78 class of catalytic priors and apply them to generalized linear
79 models. We show that a catalytic prior is proper and yields
80 stable estimates under mild conditions. Simulation studies
81 indicate the frequentist properties of the model estimator using
82 catalytic priors are comparable, and sometimes superior,
83 to existing competitive estimators. Such a prior has the
84 advantages that it is often easier to formulate and it allows for
85 simple implementation from standard software.

86 We also provide an interpretation of the catalytic prior
87 from an information theory perspective (detailed in Section 4
88 of SI).

89 **Related Priors.** The practice of using synthetic data or pseudo
90 data to define prior distributions has a long history in Bayesian
91 statistics (11). It is well-known that conjugate priors for ex-
92ponential families can be viewed as the likelihood of pseudo-
93 observations (12). Some authors have suggested formulating
94 priors by obtaining additional pseudo data from experts' knowl-
95 edge (13–15), which is not easy to use in practice when data
96 have many dimensions or when numerous models or experts
97 are being considered. Refs. (16, 17) proposed to use a conju-
98 gate Beta-distribution prior with specifically chosen values of
99 covariates to approximate a multivariate Gaussian prior for
100 the regression coefficients in a logistic regression model. A
101 complication of this approach is that the augmented dataset

102 may contain impossible values for a covariate. Another ap-
103 proach is the *expected-posterior prior* (18–20), where the prior
104 is defined as the average posterior distribution over a set of
105 imaginary data sampled from a simple predictive model. This
106 approach is designed to address the challenges in Bayesian
107 model selection. Other priors have been proposed to incor-
108 porate information from previous studies. Particularly, the
109 *power prior* (21–23) formulates an informative prior generated
110 by a power of the likelihood function of historical data. One
111 limitation of this power prior is that its properness requires
112 the covariate matrix of historical or current data to have
113 full column rank (22). Recently, the *power-expected-posterior*
114 prior was proposed to alleviate the computational challenge
115 of expected-posterior priors for model selection (24, 25). It
116 incorporates the ideas of both the expected-posterior prior and
117 the power prior, but it cannot be applied when the dimension
118 of the working model is larger than the sample size. Some
119 other priors suggested in the literature have appearances simi-
120 lar to catalytic priors. Ref. (26) propose the *reference prior*
121 that maximizes the mutual information between the data and
122 the parameter, resulting in a prior density function that looks
123 similar to that of a catalytic prior but is essentially different.
124 Ref. (27) proposed a prior based on the idea of matching
125 loss functions, which, although operationally similar to the
126 catalytic prior, is conceptually different because it requires a
127 subjective initial choice for the distribution of the data. In
128 Ref. (28), the class of *penalized complexity priors* for hierarchi-
129 cal model components is based on penalizing the complexity
130 induced by the deviation from a simpler model. The simpler
131 model there needs to be nested in the working model, which
132 is not required by the catalytic prior.

Generic Formulation of Catalytic Priors

133 **Catalytic Prior in the Absence of Covariates.** Consider the
134 data, $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$, being analyzed under a working
135 model $Y_i \stackrel{i.i.d.}{\sim} f(y \mid \theta)$ governed by unknown parameter θ .
136 Suppose a model $g(y \mid \psi)$ with unknown parameter ψ , whose
137 dimension is smaller than that of θ , is stably fitted from \mathbf{Y}
138 and results in a predictive distribution $g_*(y^* \mid \mathbf{Y})$ for future
139 data drawn from $g(y \mid \psi)$. The *synthetic-data generating dis-
140 tribution* $g_*(y^* \mid \mathbf{Y})$ is used to generate the synthetic data
141 $\{Y_i^*\}_{i=1}^M$, where M is the synthetic-sample size and the asterisk
142 superscript is used to indicate synthetic data.

143 The synthetic-data generating distribution can be specified
144 by fitting a model simpler than $f(y \mid \theta)$, but it does not
145 necessarily have to be. Examples: (1) If a Bayesian analysis
146 of the simpler model can be carried out easily, $g_*(y^* \mid \mathbf{Y})$
147 can be taken to be the posterior predictive distribution under
148 the simpler model. (2) Alternatively, one can obtain a point
149 estimate $\hat{\psi}$, and $g_*(y^* \mid \mathbf{Y}) = g(y^* \mid \hat{\psi})$ can be the plug-in
150 predictive distribution. (3) If two simpler estimated models
151 are $g_*^{(1)}(y^* \mid \mathbf{Y})$ and $g_*^{(2)}(y^* \mid \mathbf{Y})$, then $g_*(y^* \mid \mathbf{Y})$ can be
152 taken to be a mixture $w g_*^{(1)}(y^* \mid \mathbf{Y}) + (1 - w) g_*^{(2)}(y^* \mid \mathbf{Y})$ for
153 some $w \in (0, 1)$.

154 The likelihood function of θ under the working model based
155 on the synthetic data $\{Y_i^*\}_{i=1}^M$ is $\ell(\theta \mid Y^*) = \prod_{i=1}^M f(Y_i^* \mid \theta)$.
156 Because these synthetic data are not really observed data, we
157 down-weight them by raising this likelihood to a power τ/M ,
158 where $\tau > 0$ is a tuning parameter called the *prior weight*.
159 This leads to the *catalytic prior* that has an unnormalized

161 density:

$$\pi_{cat,M}(\theta | \tau) \propto \left\{ \prod_{i=1}^M f(Y_i^* | \theta) \right\}^{\tau/M}, \quad [3]$$

163 which depends on the randomly drawn synthetic data $\{Y_i^*\}_{i=1}^M$.
164 The *population catalytic prior* is formally the limit of Eq. (3)
165 as M goes to infinity:

$$\pi_{cat,\infty}(\theta | \tau) \propto \exp[\tau \mathbb{E}_{g_*} \{\log f(Y^* | \theta)\}]. \quad [4]$$

166 Here the expectation $\mathbb{E}_{g_*} \{\log f(Y^* | \theta)\}$ in Eq. (4) is taken
167 with respect to $Y^* \sim g_*(Y^* | \mathbf{Y})$. The dependence of $g_*(Y^* |$
168 $\mathbf{Y})$ on the observed \mathbf{Y} emphasizes that the catalytic prior is
169 data-dependent, like that used in Box and Cox (29) for power
170 transformations.

171 The posterior density using the catalytic prior is mathematically
172 proportional to the likelihood with both the observed
173 data and the weighted synthetic data. Thus, we can implement
174 the required Bayesian inference using standard software. For
175 instance, the maximum posterior estimate (posterior mode) is
176 the same as the MLE using the weighted augmented data and
177 can be computed by existing MLE procedures, which can be
178 a computational advantage, as illustrated in Ref. (1).

180 **Catalytic Prior with Covariates.** Let $\{(Y_i, \mathbf{X}_i)\}_{i=1}^n$ be the set of
181 n pairs of a scalar response Y_i and a p -dimensional covariate
182 vector \mathbf{X}_i ; Y_i depends on \mathbf{X}_i in the working model with
183 unknown parameter β :

$$Y_i | \mathbf{X}_i, \beta \sim f(y | \mathbf{X}_i, \beta), i = 1, 2, \dots, n. \quad [5]$$

184 Let \mathbf{Y} be the vector $(Y_1, \dots, Y_n)^\top$, and \mathbb{X} be the matrix
185 $(\mathbf{X}_1, \dots, \mathbf{X}_n)^\top$. The likelihood of these data is $f(\mathbf{Y} | \mathbb{X}, \beta) =$
186 $\prod_{i=1}^n f(Y_i | \mathbf{X}_i, \beta)$.

187 Suppose a simpler model $g(y | \mathbf{X}, \psi)$ with unknown parameter ψ is stably fitted from (\mathbf{Y}, \mathbb{X}) and results in a synthetic-data generating distribution $g_*(y | \mathbf{x}, \mathbf{Y}, \mathbb{X})$. Note that $g_*(\cdot)$ here is analogous to its use earlier except that now, in addition to the observed data, it is also conditioned on \mathbf{x} . The synthetic covariates \mathbf{X}^* will be drawn from a distribution $Q(\mathbf{x})$, which we call the *synthetic-covariate generating distribution*. We will discuss the choice of $Q(\mathbf{x})$ shortly.

188 Given the distributions $Q(\mathbf{x})$ and $g_*(y | \mathbf{x}, \mathbf{Y}, \mathbb{X})$, the catalytic prior first draws a set of synthetic data $\{(Y_i^*, \mathbf{X}_i^*)\}_{i=1}^M$ from

$$\mathbf{X}_i^* \stackrel{i.i.d.}{\sim} Q(\mathbf{x}), \quad Y_i^* | \mathbf{X}_i^* \sim g_*(y | \mathbf{X}_i^*, \mathbf{Y}, \mathbb{X}).$$

189 Hereafter we write \mathbf{Y}^* for the vector of synthetic responses
190 $(Y_1^*, \dots, Y_M^*)^\top$, and \mathbb{X}^* for the matrix of synthetic covariates
191 $(\mathbf{X}_1^*, \dots, \mathbf{X}_M^*)^\top$. The likelihood of the working model based on
192 the synthetic data $\ell(\beta | \mathbf{Y}^*, \mathbb{X}^*)$ equals $\prod_{i=1}^M f(Y_i^* | \mathbf{X}_i^*, \beta)$.
193 Because these synthetic data are not really observed, we down-weight
194 them by raising this likelihood to a power τ/M , which
195 gives the unnormalized density of the catalytic prior with
196 covariates:

$$\pi_{cat,M}(\beta | \tau) \propto \left\{ \prod_{i=1}^M f(Y_i^* | \mathbf{X}_i^*, \beta) \right\}^{\tau/M}. \quad [6]$$

197 The population catalytic prior (when $M \rightarrow \infty$) has unnormalized
198 density:

$$\pi_{cat,\infty}(\beta | \tau) \propto \exp(\tau \mathbb{E}_{Q,g_*} [\log f(Y^* | \mathbf{X}^*, \beta)]), \quad [7]$$

199 where the expectation \mathbb{E}_{Q,g_*} averages over both \mathbf{X}^* and Y^* .
200 Denote by $Z_{\tau,M}$ and $Z_{\tau,\infty}$ the integrals of the right-hand sides
201 of Eq. (6) and Eq. (7) w.r.t. β . When these integrals are
202 finite, the priors are proper, and $Z_{\tau,M}$ and $Z_{\tau,\infty}$ are their
203 normalizing constants.

204 An advantage of the catalytic prior is that the corresponding
205 posterior has the same form as the likelihood

$$\begin{aligned} \pi(\beta | \mathbb{X}, \mathbf{Y}, \tau) &\propto \pi_{cat,M}(\beta | \tau) f(\mathbf{Y} | \mathbb{X}, \beta) \\ &\propto \exp \left(\frac{\tau}{M} \sum_{i=1}^M \log(f(Y_i^* | \mathbf{X}_i^*, \beta) \right. \\ &\quad \left. + \sum_{i=1}^n \log(f(Y_i | \mathbf{X}_i, \beta) \right), \end{aligned}$$

206 which makes the posterior inference no more difficult than
207 other standard likelihood-based methods. For example, the
208 posterior mode can be easily computed as a maximum weighted
209 likelihood estimate using standard statistical software. Full
210 posterior inference can also be easily implemented by treating
211 the synthetic data as down-weighted data.

212 **Catalytic Prior for GLMs.** A generalized linear model (GLM)
213 assumes that, given a covariate vector \mathbf{X} , the response Y has
214 the following density w.r.t. some base probability measure:

$$f(y | \mathbf{X}, \beta) = \exp(t(y)\theta - b(\theta)), \quad [8]$$

215 where $t(y)$ is a sufficient statistic, and θ is the canonical
216 parameter that depends on $\eta = \mathbf{X}^\top \beta$ through $\theta = \phi(\eta)$,
217 where β is the unknown regression coefficient vector and $\phi(\cdot)$
218 is a monotone differentiable function. The mean of $t(Y)$ is
219 denoted by $\mu(\eta)$ and is equal to $b'(\phi(\eta))$.

220 When the working model is a GLM, from Eq. (7) and
221 Eq. (8), we have

$$\mathbb{E}_{Q,g_*} [\log f(Y^* | \mathbf{X}^*, \beta)]$$

$$= \mathbb{E}_Q \{ \phi(\beta^\top \mathbf{X}^*) \mathbb{E}_{g_*} [t(Y^*) | \mathbf{X}^*] - b(\phi(\beta^\top \mathbf{X}^*)) \}, \quad [9]$$

222 so that the expectation of the log-likelihood does not depend
223 on particular realizations of the synthetic response, but rather
224 on the conditional mean of the sufficient statistic under the
225 synthetic-data generating distribution. Thus, in the case of a
226 GLM (and exponential family models), instead of a specific
227 realization of the synthetic response, one only needs to use
228 the conditional mean of the sufficient statistic $\mathbb{E}_{g_*}[t(Y^*) | \mathbf{X}^*]$
229 to form a catalytic prior. This simplification reduces the
230 variability introduced by synthetic data.[†]

231 As a concrete example, consider a linear regression model
 $\mathbf{Y} = \mathbb{X}\beta + \epsilon$, where $\epsilon \sim N_n(\mathbf{0}, \sigma^2 \mathcal{I}_n)$ with known σ . Suppose
232 the synthetic-data generating model is a sub-model with the
233 estimated parameter β_0^* , and \mathbb{X}^* is the synthetic covariate
234 matrix. In this case, the catalytic prior with any positive τ
235 has a normal distribution:

$$\beta \sim N \left(\beta_0^*, \frac{\sigma^2}{\tau} \left(\frac{1}{M} (\mathbb{X}^*)^\top \mathbb{X}^* \right)^{-1} \right).$$

236 If $\lim_{M \rightarrow \infty} \frac{1}{M} (\mathbb{X}^*)^\top \mathbb{X}^* = \Sigma_{\mathbf{X}}$, the population catalytic prior
237 is

$$\beta \sim N \left(\beta_0^*, \frac{\sigma^2}{\tau} (\Sigma_{\mathbf{X}})^{-1} \right).$$

238 More details about this example can be found in SI.

239 [†]Note that in the previous example of 1970-1980 I/O code mapping, instead of the raw counts of
240 synthetic responses, their expected values $p\hat{\mu}/J$ and $p(1 - \hat{\mu})/J$ were used.

233 **Specifications of the Catalytic Prior**

234 **Generating Synthetic Covariates.** The synthetic covariate vectors are generated such that $(\mathbb{X}^*)^\top \mathbb{X}^*$ has full rank. Moreover, 235 a synthetic covariate should have the same sample space as a real covariate. The simple choice of resampling the observed covariate vectors would not guarantee the full rank of $(\mathbb{X}^*)^\top \mathbb{X}^*$; for example, if the observed covariates are rank 236 deficient, resampling would still give rank deficient $(\mathbb{X}^*)^\top \mathbb{X}^*$. 237

238 Instead, we consider one option for generating synthetic covariates: resample each coordinate of the observed covariates 239 independently. Formally, we define *the independent resampling 240 distribution* by the probability mass function

$$245 Q_0(\mathbf{x}) := \prod_j \left(\frac{1}{n} \# \{1 \leq i \leq n : (\mathbf{X}_i)_j = x_j\} \right),$$

246 for all $\mathbf{x} \in \mathcal{X}$, where \mathcal{X} is the sample space of \mathbf{X} . We use this 247 distribution for simplicity. Alternatively, if historical data are 248 available, synthetic covariates can be sampled from the historical 249 covariates. Furthermore, if some variables are naturally 250 grouped or highly correlated, one may want to resample these 251 grouped parts together. Other examples are discussed in SI.

252 **Generating Synthetic Responses.** The synthetic-data generating 253 distribution can be specified by fitting a simple model 254 $G_\Psi = \{g(y \mid \mathbf{x}, \psi) : \psi \in \Psi\}$ to the observed data. The 255 only requirement is that this simple model can be stably fit 256 by the observed data in the sense that the standard estimation 257 of ψ , using either a Bayesian or frequentist approach, 258 can lead to a well-defined predictive distribution for future 259 data. Examples include a fixed distribution and an intercept- 260 only model. G_Ψ can also be a regression model based on 261 dimension reduction, such as a principal components analysis; 262 see SI for a numerical example, which also suggests to 263 keep G_Ψ as simple as possible when the observed sample size 264 is small. For a working regression model with interactions, 265 a natural choice of G_Ψ is the sub-model with only main- 266 effects. If the main-effect model is overfitted as well, we could 267 use a mixed synthetic-data generating distribution, such as 268 $g_*(y \mid \mathbf{x}, \mathbf{Y}, \mathbb{X}) = 0.5 g_{*,1}(y \mid \mathbf{x}, \mathbf{Y}, \mathbb{X}) + 0.5 g_{*,0}(y \mid \mathbf{x}, \mathbf{Y}, \mathbb{X})$, 269 where $g_{*,1}$ and $g_{*,0}$ are the predictive distributions of the pre- 270 liminarily fitted main-effect model and intercept-only model, 271 respectively. G_Ψ can also be chosen using additional knowl- 272 edge, such as a sub-model that includes a few important 273 covariates that have been identified in previous studies, or if 274 domain experts have opinions on the range of possible values 275 of certain model parameters, then the parameter space Ψ can 276 be constrained accordingly.

277 Sometimes it is beneficial to draw multiple synthetic re- 278 sponses for each sampled synthetic covariate vector. We name 279 this sampling the *stratified synthetic data generation*. It could 280 help reduce variability introduced by synthetic data.

281 **Sample Size of Synthetic Data.** Theorem 5 below quantifies 282 how fast the randomness in the catalytic prior diminishes as 283 the synthetic-sample size M increases. One implication is that 284 for linear regression with binary covariates, if $M \geq \frac{4p^3}{\epsilon^2} \log(\frac{p}{\delta})$, 285 then the Kullback-Leibler divergence between the catalytic 286 prior $\pi_{cat,M}$ and its limit $\pi_{cat,\infty}$ is at most ϵ with probability 287 at least $1 - \delta$. Such a bound can help choose the magnitude of 288 M . When the prior needs to be proper, we suggest taking M

289 larger than 4 times the dimension of β (based on Theorem 1 290 and Proposition 2 below).

291 **Weight of Synthetic Data.** The prior weight τ controls how 292 much the posterior inference relies on the synthetic data be- 293 cause it can be interpreted as the effective prior sample size. 294 Here we provide two guidelines for systematic specifications 295 of τ .

296 **Frequentist Predictive Risk Estimation.** Choose a value of τ using 297 the following steps: (1) Compute the posterior mode $\hat{\beta}(\tau)$ 298 for various values of τ . (2) Choose a discrepancy function 299 $D(y_0, \hat{\mu})$ that measures how well a prediction $\hat{\mu}$ predicts a 300 future response y_0 . (3) Find an appropriate criterion function 301 $\Lambda(\tau)$ that estimates the expected (in-sample) prediction error, 302 for a future response Y_0 based on $\hat{\beta}(\tau)$, and (4) Pick the value 303 of τ that minimizes $\Lambda(\tau)$. See Section 2.C.1 in SI for a detailed 304 discussion.

305 The discrepancy $D(y_0, \hat{\mu})$ measures the error of a prediction 306 $\hat{\mu}$ for a future response Y_0 that takes value y_0 . We consider 307 here discrepancy functions of the form

$$308 D(y_0, \hat{\mu}) := a(\hat{\mu}) - \lambda(\hat{\mu})y_0 + c(y_0), \quad [10]$$

309 and define $\mathbf{D}(Y_0, \hat{\mu}) := \frac{1}{n} \sum_{i=1}^n D(Y_{0,i}, \hat{\mu}_i)$. This class is 310 general enough to include squared error, classification error and 311 deviance for GLMs: (a) squared error: $D(y_0, \hat{\mu}) = (y_0 - \hat{\mu})^2 = \hat{\mu}^2 - 2y_0\hat{\mu} + y_0^2$; (b) classification error: $D(y_0, \hat{\mu}) = \mathbf{1}_{y_0 \neq \hat{\mu}} = \hat{\mu} - 2y_0\hat{\mu} + y_0$ for any y_0 and $\hat{\mu}$ in $\{0, 1\}$; (c) deviance for GLMs: 312 $D(y_0, \hat{\mu}) = b(\hat{\theta}) - y_0\hat{\theta} + \sup_{\theta} (y_0\theta - b(\theta))$, where $\hat{\theta} = (b')^{-1}(\hat{\mu})$. 313

314 The criterion function $\Lambda(\tau)$ is an estimate of the expectation 315 of the (in-sample) prediction error. Such an estimate can be 316 obtained by using the parametric bootstrap. Take a bootstrap 317 sample of the response vector \mathbf{Y}^{boot} from the distribution 318 $f(\mathbf{y} \mid \mathbb{X}, \hat{\beta}^0)$, where $\hat{\beta}^0 = \hat{\beta}(\tau_0)$ is a preliminary estimate, and 319 denote by $\hat{\beta}^{boot}(\tau)$ the posterior mode based on data $(\mathbf{Y}^{boot}, \mathbb{X})$ 320 with the catalytic prior. The bootstrap criterion function is 321 given by

$$322 \Lambda(\tau) = \mathbf{D}(\mathbf{Y}, \hat{\mu}_\tau) + \frac{1}{n} \sum_{i=1}^n \text{Cov}(\lambda(\hat{\mu}_{\tau,i}^{boot}), Y_i^{boot}), \quad [11]$$

323 where $\hat{\mu}_{\tau,i} = \mu(\mathbf{X}_i^\top \hat{\beta}(\tau))$ and $\hat{\mu}_{\tau,i}^{boot} = \mu(\mathbf{X}_i^\top \hat{\beta}^{boot}(\tau))$. 324 See SI for a detailed derivation. In practice, the term 325 $\text{Cov}(\lambda(\hat{\mu}_{\tau,i}^{boot}), Y_i^{boot})$ is numerically computed by sampling 326 \mathbf{Y}^{boot} repeatedly. Based on our experiments with linear and 327 logistic models, the default choices of the initial values can be 328 $\tau_0 = 1$ for linear regression and $\tau_0 = p/4$ for other cases. See 329 SI for a mathematical argument.

330 The costly bootstrap repetition step to numerically compute 331 $\text{Cov}(\lambda(\hat{\mu}_{\tau,i}^{boot}), Y_i^{boot})$ can be avoided in two special cases (see 332 SI for more discussion):

- 333 1. If Y_i follows a normal distribution and $\lambda(\hat{\mu}_{\tau,i})$ is smooth 334 in y_i , then the *Stein's unbiased risk estimate* yields

$$335 \Lambda(\tau) = \mathbf{D}(\mathbf{Y}, \hat{\mu}_\tau) + \frac{1}{n} \sum_{i=1}^n \text{Var}(Y_i) \mathbb{E} \frac{\partial \lambda(\hat{\mu}_{\tau,i})}{\partial y_i}. \quad [12]$$

336 In particular, when square error is considered and if $\hat{\mu}_\tau$ 337 can be written as $\hat{\mu}_\tau = \mathbf{H}_\tau \cdot \mathbf{Y} + \mathbf{c}_\tau$, the risk estimate is

$$338 \Lambda(\tau) = \|\mathbf{Y} - \hat{\mu}_\tau\|^2 + \frac{2}{n} \sum_{i=1}^n \text{Var}(Y_i) \mathbf{H}_\tau(i, i). \quad [13]$$

- 325 2. When responses are binary, say 0 or 1, let $\mathbf{Y}^{\triangle i}$ be a copy
 326 of \mathbf{Y} but with Y_i replaced by $1 - Y_i$, and let $\hat{\beta}^{\triangle i}(\tau)$ be the
 327 posterior mode based on data $(\mathbf{X}, \mathbf{Y}^{\triangle i})$ with the catalytic
 328 prior. The *Steinian* estimate (30) is given by

$$329 \mathbf{D}(\mathbf{Y}, \hat{\mu}_\tau) + \frac{1}{n} \sum_{i=1}^n \hat{\mu}_i^0 (1 - \hat{\mu}_i^0) (2Y_i - 1) (\lambda(\hat{\mu}_{\tau,i}) - \lambda(\hat{\mu}_{\tau,i}^{\triangle i})) , \quad [14]$$

330 where $\hat{\mu}_i^0 = \mu(\mathbf{X}_i^\top \hat{\beta}^0)$, and $\hat{\mu}_{\tau,i}^{\triangle i} = \mu(\mathbf{X}_i^\top \hat{\beta}^{\triangle i}(\tau))$.

331 **Bayesian Hyperpriors.** An alternative way to specify the prior
 332 weight τ is to consider a joint catalytic prior for (τ, β) :

$$333 \pi_{\alpha, \gamma}(\tau, \beta) \propto \Gamma_{\alpha, \gamma}(\tau) \left\{ \prod_{i=1}^M f(Y_i^* | \mathbf{X}_i^*, \beta) \right\}^{\tau/M} , \quad [15]$$

334 where $\Gamma_{\alpha, \gamma}(\tau)$ is a function defined as follows for positive scalar
 335 hyperparameters α and γ . Denote

$$336 \kappa := \sup_{\beta \in \mathbb{R}^p} \frac{1}{M} \sum_{i=1}^M \log f(Y_i^* | \mathbf{X}_i^*, \beta).$$

337 For linear regression, the function $\Gamma_{\alpha, \gamma}(\tau)$ can be taken to be

$$338 \Gamma_{\alpha, \gamma}(\tau) = \tau^{\frac{p+\alpha}{2}-1} e^{-\tau(\kappa+\gamma^{-1})}. \quad [16]$$

339 and for other models,

$$340 \Gamma_{\alpha, \gamma}(\tau) = \tau^{p+\alpha-1} e^{-\tau(\kappa+\gamma^{-1})}. \quad [17]$$

341 The form of $\Gamma_{\alpha, \gamma}(\tau)$ is chosen mainly for practical convenience;
 342 by separating the dependence on p and κ , we have
 343 meaningful interpretations for α and γ . For GLMs, prior
 344 moments of β up to order α exist, and γ controls the exponential
 345 decay of the prior density of τ ; see Theorem 4. For linear
 346 regression, the marginal prior for β induced by Eq. (15) is a
 347 multivariate t-distribution centered around the MLE for the
 348 synthetic data with covariance matrix $\frac{2\sigma^2}{\alpha\gamma} \cdot (\frac{1}{M}(\mathbf{X}^*)^\top \mathbf{X}^*)^{-1}$
 349 and degrees of freedom α . The analysis in Theorem 4 reveals
 350 how the parameters α and γ affect the joint prior. Roughly
 351 speaking, a larger value of α (or γ) tends to pull the working
 352 model more towards the simpler model. Admittedly, it appears
 353 impossible to have a single choice that works the best in all
 354 scenarios. We recommend $(\alpha, \gamma) = (2, 1)$ as a simple default
 355 choice based on our numerical experiments.

356 Illustration of Methods

357 **Logistic Regression.** We illustrate the catalytic prior using
 358 logistic regression. Another example using linear regression
 359 is presented in SI. Here the mean of Y depends on the linear
 360 predictor $\eta = \mathbf{X}^\top \beta$ through $\mu = e^\eta / (1 + e^\eta)$. Suppose the
 361 synthetic-data generating model includes only the intercept,
 362 so it is Bernoulli(μ_0), where a simple estimate of μ_0 is given
 363 by $\hat{\mu}_0 = (1/2 + \sum_{i \leq n} Y_i) / (1 + n)$. The synthetic response
 364 vector \mathbf{Y}^* can be taken to be $\hat{\mu}_0 \cdot \mathbf{1}_M$, and each synthetic
 365 covariate vector \mathbf{X}_i^* is drawn from the independent resampling
 366 distribution; this prior is proper when $(\mathbf{X}^*)^\top \mathbf{X}^*$ is positive
 367 definite according to Theorem 1.

Numerical Example. We first generate the observed covariates \mathbf{X}_i
 368 by drawing a Gaussian random vector \mathbf{Z}_i whose components
 369 have mean 0, variance 1 and common correlation $\rho = 0.5$; set
 370

$$371 \mathbf{X}_{i,j} = \begin{cases} 2 \cdot \mathbf{1}_{\mathbf{Z}_{i,j} > 0} - 1, & 2j < p \\ \mathbf{Z}_{i,j}, & 2j \geq p. \end{cases}$$

372 This process yields covariate vectors that have dependent components
 373 and have both continuous and discrete components
 374 as one would encounter in practical logistic regression problems.
 375 We consider three different sparsity levels and three
 376 different amplitudes of the regression coefficient β in the
 377 underlying model. More precisely, β is specified through scaling
 378 an initial coefficient $\beta^{(0)}$ that accommodates different levels
 379 of sparsity. Each coordinate of $\beta^{(0)}$ is either 1 or 0. ζ proportion
 380 of the coordinates of $\beta^{(0)}$ are randomly selected and
 381 set to 1, and the remaining $1 - \zeta$ proportion are set to 0,
 382 where ζ is the level of *non-sparsity* and is set at 1/4, 1/2,
 383 3/4. This factor controls how many covariates actually affect
 384 the response. Then the amplitude of β is specified indirectly:
 385 $\beta_0 = c_1$, $\beta_{1:(p-1)} = c_2 \beta_{1:(p-1)}^{(0)}$, where parameters (c_1, c_2) are
 386 chosen such that the oracle classification error r (the expected
 387 classification error of the classifier given by the true β) is equal
 388 to 0.1, 0.2, 0.3. Here $r = \mathbb{E}_{\mathbf{X}} (\min(\mathbf{P}_\beta(Y=1), \mathbf{P}_\beta(Y=0))) =$
 389 $\mathbb{E}_{\mathbf{X}} (1 + \exp(|\mathbf{X}^\top \beta|))^{-1}$ is numerically computed by sampling
 390 2000 extra covariate vectors. The value of r represents how
 391 far apart the class $Y = 1$ is from the class $Y = 0$, and small
 392 values of r correspond to large amplitudes of β .

393 In this example, the number of covariates is 16, so the dimension
 394 of β is $p = 17$, and the sample size is $n = 30$. We use the
 395 predictive binomial deviance, $\mathbb{E}_{\mathbf{X}_0} [D(\mu(\mathbf{X}_0^\top \beta), \mu(\mathbf{X}_0^\top \hat{\beta}))]$,
 396 where $D(a, b) = a \log(a/b) + (1 - a) \log((1 - a)/(1 - b))$ measures the discrepancy
 397 between two Bernoulli distributions with probability a and b respectively, to evaluate the predictive
 398 performance of $\hat{\beta}$. The expectation $\mathbb{E}_{\mathbf{X}_0}$ is computed by
 399 sampling 1000 extra independent copies of \mathbf{X}_0 from the same
 400 distribution that generates the observed covariates.

401 To specify catalytic priors, we use the generating distributions
 402 for synthetic data just described, and fix M at 400. The first
 403 estimator of β is the posterior mode of β with $\tau = \hat{\tau}_{boot}$
 404 selected by predictive risk estimation via the bootstrap with
 405 deviance discrepancy (denoted as *Cat.Boot.*). This estimator
 406 can be computed as the MLE with the weighted augmented
 407 data. The second estimator of β is the coordinate-wise posterior
 408 median of β with the joint prior $\pi_{\alpha=2, \gamma=1}$ (denoted as
 409 *Cat.Joint*). The posterior median is used here because there is
 410 no guarantee that the posterior distribution of β is uni-modal
 411 in this case. These estimators are compared with two alternatives:
 412 the MLE and the posterior mode with the Cauchy prior
 413 (31) (calculated by the authors' R package *bayesglm*).

414 Table 1 presents the average predictive binomial deviance
 415 over 1600 simulations in each cell. The column *Comp.Sep.*
 416 shows how often complete separation occurs in the datasets;
 417 when complete separation occurs, the MLE does not exist but
 418 a pseudo-MLE can be algorithmically computed if the change
 419 in the estimate is smaller than 10^{-8} within 25 iterations; the
 420 column of MLE averages across only the cases where either
 421 MLE or pseudo-MLE exists. In Table 1, the boldface corresponds
 422 to the best performing method under each simulation scenario.
 423 Based on this table, the catalytic prior with $\hat{\tau}_{boot}$
 424 predicts the best and the MLE predicts the worst in all cases
 425 considered. Although the Cauchy prior seems to perform close
 426

Setting		Comp. Sep.	Performance of Methods				
ζ	r		Cat. Boot.	Cat. Joint	Cauchy	MLE (pseudo)	
1/4	0.1	100%	Mean $SE \times 10^3$	1.692 (6.8)	1.772 (6.7)	1.793 (6.7)	2.081 (8.7)
	0.2	98%	Mean $SE \times 10^3$	0.675 (5.2)	0.769 (5.0)	0.802 (5.0)	1.123 (7.2)
0.3	91%	Mean $SE \times 10^3$	0.297 (2.3)	0.399 (2.0)	0.445 (1.9)	0.751 (7.3)	
	0.1	100%	Mean $SE \times 10^3$	1.661 (3.9)	1.742 (3.8)	1.749 (3.8)	2.048 (5.0)
2/4	0.2	98%	Mean $SE \times 10^3$	0.648 (2.5)	0.743 (2.2)	0.771 (2.0)	1.107 (3.4)
	0.3	92%	Mean $SE \times 10^3$	0.287 (2.1)	0.392 (1.8)	0.438 (1.7)	0.748 (7.1)
3/4	0.1	100%	Mean $SE \times 10^3$	1.664 (4.0)	1.746 (3.9)	1.749 (3.8)	2.052 (4.9)
	0.2	99%	Mean $SE \times 10^3$	0.649 (2.5)	0.745 (2.2)	0.771 (2.0)	1.104 (3.4)
	0.3	91%	Mean $SE \times 10^3$	0.287 (2.1)	0.391 (1.9)	0.435 (1.7)	0.738 (7.3)

Table 1. Mean and standard error of predictive binomial deviance of different methods. The first two columns are the settings of the simulation: ζ is the non-sparsity, and r is the oracle prediction error. The column of Comp.Sep. shows how often complete separation occurs in the datasets. The last four columns report the mean and standard error of the predictive binomial deviance of the different methods, which are the catalytic posterior mode with $\hat{\tau}_{boot}$, denoted by Cat.Boot., the posterior median under joint catalytic prior, denoted by Cat.Joint, the Cauchy posterior mode, denoted by Cauchy, and the MLE. The boldface corresponds to the best performing method in each simulation scenario.

Difference between the error of Cauchy and that of Cat.Joint		
ζ	r	Mean $SE \times 10^3$
1/4	0.1	0.021 0.98
	0.2	0.033 0.91
	0.3	0.047 0.86
1/2	0.1	0.007 0.79
	0.2	0.028 0.85
	0.3	0.046 0.84
3/4	0.1	0.003 0.76
	0.2	0.026 0.83
	0.3	0.044 0.82

Table 2. Mean and standard error of the difference in predictive binomial deviance between the Cauchy posterior mode and the joint catalytic posterior median. ζ is the non-sparsity; r is the oracle prediction error.

Setting	ζ	r	Performance of Methods		
			Cat.Boot	Cat.Joint	Cauchy
1/4	0.1	Cover	90.5%	88.1%	90.1%
		Width	3.5	2.9	3.3
		Cover	93.3%	97.2%	98.0%
	0.2	Width	2.8	2.7	3.0
		Cover	95.0%	97.6%	97.6%
		Width	2.2	2.4	2.8
2/4	0.1	Cover	89.8%	85.7%	86.2%
		Width	3.5	2.9	3.2
		Cover	93.4%	97.5%	98.4%
	0.2	Width	2.7	2.7	3.0
		Cover	95.7%	97.7%	97.7%
		Width	2.1	2.4	2.8
3/4	0.1	Cover	89.4%	85.6%	86.1%
		Width	3.5	2.9	3.2
		Cover	93.9%	97.6%	98.6%
	0.2	Width	2.7	2.7	3.0
		Cover	95.9%	97.8%	97.8%
		Width	2.1	2.4	2.7

Table 3. Average coverage probability (%) and width of 95% posterior intervals under the catalytic prior with $\hat{\tau}_{boot}$, the joint catalytic prior, and Cauchy prior. ζ is the non-sparsity; r is the oracle prediction error.

Table 3 presents the average coverage probabilities (in percentage) and widths of the 95% nominal intervals for β_j averaging over j . Because all the intervals given by the MLE have widths too large to be useful (thousands of times wider than those given by the other methods), we do not report them in this table. The intervals from the other three priors are reasonably short in all cases and have coverage rates not far from the nominal levels. Specifically, the intervals given by the Cauchy prior and the joint catalytic prior tend to over-cover when the true β has small amplitudes ($r = 0.2$ or 0.3) and tend to under-cover when β has large amplitudes ($r = 0.1$), whereas the intervals given by the catalytic prior with $\hat{\tau}_{boot}$ perform more consistently. This example, together with more results given in SI, illustrates that, for logistic regression, the catalytic prior is at least as good as the Cauchy prior. The SI also illustrates the performance of the catalytic prior in linear regression, where it is at least as good as ridge regression. Catalytic priors thus appear to provide a general framework for prior construction over a broad range of models.

Theoretical Properties of Catalytic Priors

We show the properness and the convergence of a catalytic prior when the working model is a GLM. Without loss of generality, we assume the sufficient statistic in the GLM formula Eq. (8) is $t(y) = y$; otherwise, we can let the response be $Y' = t(Y)$ and proceed. We assume that every covariate has at least two different observed values. Denote by \mathcal{Y} the nonempty interior of the convex hull of the support of the model density in Eq. (8). Our results apply to any positive prior weight τ .

Properness. A proper prior is needed for many Bayesian inferences, such as model comparison using Bayes factors (32). We show that catalytic priors, population catalytic priors, and joint catalytic priors are generally proper, with proofs in SI.

Theorem 1. Suppose (1) $\phi(\cdot)$ satisfies $\inf_{\eta \neq 0} |\phi(\eta)|/\eta > 0$, (2) the synthetic covariate matrix \mathbb{X}^* has full column rank,

to the joint catalytic prior, Table 2 shows that the prediction based on the joint catalytic prior is statistically significantly better than that of the Cauchy prior (Table 2 directly calculates the difference of the prediction errors between the Cauchy prior and the joint catalytic prior and shows that the difference is significantly positive with Bonferroni-corrected p-value smaller than 0.02). Tables 1 and 2 focus on predictive binomial deviance. Section 3.D in SI considers other error measurements, including the classification error and the AUC (Area Under Curve), where a similar conclusion can be drawn regarding the performance of different methods: predictions based on catalytic priors are generally much better than those based on the MLE and are often better than those based on the Cauchy prior.

475 and (3) each synthetic response Y_i^* lies in \mathcal{Y} or there exists a
 476 linearly independent subset $\{X_{i_k}^*\}_{k=1}^p$ of the synthetic covariate
 477 vectors such that the average of synthetic responses with the
 478 same $X_{i_k}^*$ lies in \mathcal{Y} . Then the catalytic prior is proper for any
 479 $\tau > 0$.

480 The condition $\inf_{\eta \neq 0} |\phi(\eta)/\eta| > 0$ is satisfied for the canonical
 481 link for any GLM, and also for the commonly used probit
 482 link and the complementary log-log link in binary regression.
 483 The condition that \mathbb{X}^* has full column rank holds with high
 484 probability according to the following result.

485 **Proposition 2.** *If each synthetic covariate vector is drawn
 486 from the independent resampling distribution, then there exists
 487 a constant $c > 0$ that only depends on the observed \mathbb{X} such that
 488 for any $M > p$, with probability at least $1 - 2\exp(-cM)$, the
 489 synthetic covariate matrix \mathbb{X}^* has full column rank.*

490 Population catalytic priors are also proper.

491 **Theorem 3.** *Suppose (1) $\phi(\cdot)$ satisfies $\inf_{\eta \neq 0} |\phi(\eta)/\eta| > 0$,
 492 (2) the synthetic covariate vector is drawn from the independent
 493 resampling distribution, and (3) there exists a compact subset
 494 $\mathcal{Y}^{com} \subset \mathcal{Y}$ such that $\mathbf{P}(Y^* \in \mathcal{Y}^{com}) = 1$. Then the population
 495 catalytic prior is proper for any $\tau > 0$.*

496 The following result shows the properness of the joint prior
 497 $\pi_{\alpha, \gamma}(\tau, \beta)$ in Eq. (15) and the role of the hyperparameters.

498 **Theorem 4.** *Suppose α and γ are positive. If $\Gamma_{\alpha, \gamma}(\tau)$ equals
 499 Eq. (16) for linear regression or equals Eq. (17) for other
 500 generalized linear models. Then under the same condition
 501 as Theorem 1, (1) the joint prior is proper; (2) for any $m \in$
 502 $(0, \alpha)$, the m^{th} moment of β exists; (3) $\lim_{\tau \rightarrow \infty} \frac{1}{\tau} \log h_{\alpha, \gamma}(\tau) =$
 503 $-1/\gamma < 0$, where $h_{\alpha, \gamma}(\tau)$ denotes the marginal prior on τ .*

504 **Convergence to the Population Catalytic Prior.** When
 505 synthetic-sample size, M , is large enough, the randomness in
 506 the synthetic data will not affect the catalytic prior regardless
 507 of the observed real sample size because, as a distribution of
 508 the parameters, the catalytic prior converges to the population
 509 catalytic prior.

510 We can quantify how fast the catalytic prior, as a random
 511 distribution, converges to the population catalytic prior by
 512 establishing an explicit upper bound on the distance between
 513 these two distributions in terms of M . This result shows how
 514 large M needs to be so that the randomness in the synthetic
 515 data no longer influentially change the prior. We present here
 516 a simplified version of the theoretical result; for precise and
 517 detailed statements, see SI.

518 **Theorem 5.** *Under mild regularity conditions,*

1. *For any given τ and p , there exists a constant C_1 , such that for any small positive ϵ_0 , ϵ_1 , and any $M \geq C_1 \left(1 + \log^2\left(\frac{1}{\epsilon_1}\right)\right) \frac{1}{\epsilon_1^2} \log\left(\frac{1}{\epsilon_0}\right)$, with probability at least $1 - \epsilon_0$ the total variation distance between the catalytic prior and the population catalytic prior is bounded by*

$$d_{TV}(\pi_{cat, \infty}, \pi_{cat, M}) \leq \epsilon_1.$$

2. *If the working model is linear regression with Gaussian noise, then there exists a constant C_2 that only depends on the observed covariates, such that for any $\epsilon_0 > 0$ and any $M > \frac{16}{9} C_2^2 p \log\left(\frac{p}{\epsilon_0}\right)$, with probability at least $1 - \epsilon_0$,*

523 *the KL divergence between the catalytic prior and the
 524 population catalytic prior with any $\tau > 0$ is bounded by*

$$KL(\pi_{cat, \infty}, \pi_{cat, M}) \leq 2C_2 \sqrt{\frac{1}{M} p^3 \log\left(\frac{p}{\epsilon_0}\right)}. \quad 525$$

526 **Data Availability.** All the data used in the article are simula-
 527 tion data. The details, including the models to generate the
 528 simulation data, are described in the *Illustration of Methods*
 529 section and the *Additional Simulations* section of the SI.

530 Discussion

531 The class of catalytic prior distributions stabilizes the estima-
 532 tion of a relatively complicated working model by augmenting
 533 the actual data with synthetic data drawn from the predictive
 534 distribution of a simpler model (including but not limited to
 535 a sub-model of the working model). Our theoretical work
 536 and simulation-based evidence suggest that the resulting in-
 537 ferences using standard software, which treat the augmented
 538 data just like actual data, have competitive and sometimes
 539 clearly superior frequency operating characteristics, compared
 540 to inferences based on alternatives that have been previously
 541 proposed. Moreover, catalytic priors are generally easier to
 542 formulate because they are based on hypothetical smoothed
 543 data that resemble the actual data. Two tuning constants,
 544 M and τ , require selection, and wise choices for them appear
 545 to be somewhat model dependent, for example, differing for
 546 linear and logistic regressions, both of which are considered
 547 here. We anticipate that catalytic priors will find broad applica-
 548 tion, especially as more complex Bayesian models are fit to
 549 more and more complicated datasets. Some open questions
 550 for future investigation include: (1) how to apply the catalytic
 551 priors to model selection, (2) how to study the asymptotic
 552 properties when both the sample size and the dimension of
 553 the working model go to infinity — in such a regime, it is also
 554 interesting to investigate what the simple model should be in
 555 order to achieve good bias-variance tradeoffs.

556 **ACKNOWLEDGMENTS.** D.B.R.'s research is supported in part
 557 by National Science Foundation Grant IIS-1409177, National Insti-
 558 tutes of Health Grant 1R01AI140854, and Office of Naval Research
 559 Grant N00014-17-1-2131. S.C.K.'s research is supported in part by
 560 National Science Foundation Grant DMS-1810914.

1. Clogg CC, Rubin DB, Schenker N, Schulte B, Weidman L (1991) Multiple imputation of industry and occupation codes in census public-use samples using Bayesian logistic regression. *Journal of the American Statistical Association* 86(413):68–78.
2. Rubin DB (2004) *Multiple imputation for nonresponse in surveys*. (John Wiley & Sons) Vol. 81.
3. Hirano K, Imbens GW, Rubin DB, Zhou XH (2000) Assessing the effect of an influenza vaccine in an encouragement design. *Biostatistics* 1(1):69–88.
4. Day NE, Kerridge DF (1967) A general maximum likelihood discriminant. *Biometrics* pp. 313–323.
5. Albert A, Anderson JA (1984) On the existence of maximum likelihood estimates in logistic regression models. *Biometrika* 71(1):1–10.
6. Firth D (1993) Bias reduction of maximum likelihood estimates. *Biometrika* 80(1):27–38.
7. Heinze G, Schemper M (2002) A solution to the problem of separation in logistic regression. *Statistics in Medicine* 21(16):2409–2419.
8. Rubin DB, Schenker N (1987) Logit-based interval estimation for binomial data using the Jeffrey's prior. *Sociological Methodology* 17:131–144.
9. Chen MH, Ibrahim JG, Kim S (2008) Properties and implementation of Jeffrey's prior in binomial regression models. *Journal of the American Statistical Association* 103(484):1659–1664.
10. Goodman LA (1968) The analysis of cross-classified data: Independence, quasi-independence, and interactions in contingency tables with or without missing entries: Ra fisher memorial lecture. *Journal of the American Statistical Association* 63(324):1091–1131.
11. Good IJ (1983) *Good thinking: The foundations of probability and its applications*. (U of Minnesota Press).
12. Raiffa H, Schlaifer R (1961) *Applied statistical decision theory*. (Wiley Cambridge).
13. Kadane JB, Dickey JM, Winkler RL, Smith WS, Peters SC (1980) Interactive elicitation of opinion for a normal linear model. *Journal of the American Statistical Association* 75(372):845–854.

- 587 14. Bedrick EJ, Christensen R, Johnson W (1996) A new perspective on priors for generalized
588 linear models. *Journal of the American Statistical Association* 91(436):1450–1460.
589 15. Bedrick EJ, Christensen R, Johnson W (1997) Bayesian binomial regression: Predicting sur-
590 vival at a trauma center. *The American Statistician* 51(3):211–218.
591 16. Greenland S, Christensen R (2001) Data augmentation priors for Bayesian and semi-Bayes
592 analyses of conditional-logistic and proportional-hazards regression. *Statistics in Medicine*
593 20(16):2421–2428.
594 17. Greenland S (2001) Putting background information about relative risks into conjugate prior
595 distributions. *Biometrics* 57(3):663–670.
596 18. Iwaki K (1997) Posterior expected marginal likelihood for testing hypotheses. *Journal of Eco-
597 nomics, Asia University* 21:105–134.
598 19. Neal RM (2001) Transferring prior information between models using imaginary data, (De-
599 partment of Statistics, University of Toronto), Technical Report 0108.
600 20. Pérez JM, Berger JO (2002) Expected-posterior prior distributions for model selection.
601 *Biometrika* 89(3):491–512.
602 21. Ibrahim JG, Chen MH (2000) Power prior distributions for regression models. *Statistical
603 Science* 15(1):46–60.
604 22. Chen MH, Ibrahim JG, Shao QM (2000) Power prior distributions for generalized linear mod-
605 els. *Journal of Statistical Planning and Inference* 84(1-2):121–137.
606 23. Ibrahim JG, Chen MH, Sinha D (2003) On optimality properties of the power prior. *Journal of
607 the American Statistical Association* 98(461):204–213.
608 24. Fouskakis D, Ntzoufras I, Draper D (2015) Power-expected-posterior priors for variable selec-
609 tion in Gaussian linear models. *Bayesian Analysis* 10(1):75–107.
610 25. Fouskakis D, Ntzoufras I, Perrakis K (2018) Power-expected-posterior priors for generalized
611 linear models. *Bayesian Analysis* 13(3):721–748.
612 26. Bernardo JM (1979) Reference posterior distributions for Bayesian inference. *Journal of the
613 Royal Statistical Society: Series B (Methodological)* 41(2):113–128.
614 27. Brown PJ, Walker SG (2012) Bayesian priors from loss matching. *International Statistical
615 Review* 80(1):60–82.
616 28. Simpson D, Rue H, Riebler A, Martins TG, Sørbye SH (2017) Penalising model component
617 complexity: A principled, practical approach to constructing priors. *Statistical science* 32(1):1–
618 28.
619 29. Box GE, Cox DR (1964) An analysis of transformations. *Journal of the Royal Statistical
620 Society: Series B (Methodological)* 26(2):211–243.
621 30. Efron B (2004) The estimation of prediction error: covariance penalties and cross-validation.
622 *Journal of the American Statistical Association* 99(467):619–632.
623 31. Gelman A, Jakulin A, Pittau MG, Su YS (2008) A weakly informative default prior distribution
624 for logistic and other regression models. *The Annals of Applied Statistics* 2(4):1360–1383.
625 32. Kass RE, Raftery AE (1995) Bayes factors. *Journal of the American Statistical Association*
626 90(430):773–795.