

---

# DegenDetector: Symbolic Recovery of Parameter Degeneracies in Bayesian Posteriors

---

Chaipat Tirapongprasert<sup>1\*</sup> Matthew Ho<sup>1</sup>

<sup>1</sup>Department of Astronomy, Columbia University, New York, NY, 10027

## Abstract

We introduce DEGENDETECTOR, a framework for identifying and characterizing parameter degeneracies in posterior distributions as closed-form symbolic equations. By combining mutual information screening with alternating symbolic regression, we facilitate automated and interpretable identification of degenerate relationships without domain-specific input. While standard tools such as corner plots can indicate that correlations exist, they do not reveal the underlying functional form; DEGENDETECTOR fills this gap by expressing multi-parameter degeneracies as closed-form equations, providing interpretable structure that scales to high-order parameter spaces.

## 1 Introduction

A central goal in physics is to constrain the parameters of physical models using experimental data, which, within Bayesian statistics, amounts to computing a posterior distribution over model parameters. Markov Chain Monte Carlo [MCMC; Goodman and Weare, 2010a, Skilling, 2006] algorithms allow for obtaining such posteriors provided that we can write down the analytic likelihood. Complex physical phenomena, however, often necessitate numerically intensive simulations to evaluate the likelihood. By training a neural surrogate to mimic simulation outputs, simulation-based inference [SBI; Papamakarios et al., 2019] has become the standard framework to address these intractable likelihoods.

Parameter degeneracy is a persistent challenge in both settings. This occurs when the observed data constrains a combination of parameters more tightly than the individual parameters themselves, causing posterior samples to concentrate along a lower-dimensional manifold. Some canonical cosmological examples include the CMB power spectrum [Efstathiou and Bond, 1999] and weak-lensing degeneracies [Jain and Seljak, 1997]. Left undetected, such degeneracies have serious downstream consequences, such as producing deceptively narrow error bars for coupled parameter combinations that explode for individual parameters [Efstathiou and Bond, 1999]. MCMC marginalization over degenerate directions can also produce one-dimensional posteriors that misrepresent the best-fit region entirely, yielding confidence intervals that are artifacts of the projection rather than reflections of the data themselves [Colgáin et al., 2025].

Detecting degeneracies has proved to be a difficult task [Jasche and Wandelt, 2013, Fluri et al., 2021], especially for a joint distribution of more than two parameters, as these high-order degeneracies do not manifest in the pairwise projections given by corner plots and other standard diagnostic tools. To address these challenges, we present DEGENDETECTOR, an automated pipeline that (a) identifies which subsets of parameters contribute to the degeneracy and (b) expresses such degeneracy as a closed-form symbolic equation that can be used to inform reparameterization [Villa et al., 2025] and decorrelate the degeneracy. We demonstrate our algorithm on synthetic posteriors with known analytical degeneracies and apply it to posterior chains from Planck CMB measurements [Aghanim et al., 2020]. All code is publicly available on Github.<sup>1</sup>

---

<sup>1</sup>[https://github.com/chaipattira/deggen\\_detector](https://github.com/chaipattira/deggen_detector)

## 2 Methods

Consider a Bayesian inference problem where we have used MCMC or neural SBI methods to generate  $N$  samples  $\{\boldsymbol{\theta}^{(i)}\}_{i=1}^N$  from a posterior distribution  $p(\boldsymbol{\theta} \mid \mathbf{d})$  over an  $M$ -dimensional parameter space  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_M)$ .

**Degeneracies** The samples are said to exhibit  $k$ -parameter degeneracy if there exists a smooth function  $F(\theta_{j_1}, \dots, \theta_{j_k})$  such that most posterior mass lies within an  $\varepsilon$ -neighborhood of

$$\mathcal{M}_c = \{ \boldsymbol{\theta} \in \mathbb{R}^M : F(\theta_{j_1}, \dots, \theta_{j_k}) = c \}, \quad (1)$$

under some metric on parameter space, where  $(j_1, \dots, j_k) \subseteq [1, M]$  is a tuple of the indices of the  $k$  parameters involved in the degeneracy. Generally, the level set  $\mathcal{M}_c$  is a  $(k - 1)$ -dimensional submanifold of  $\mathbb{R}^M$ , rendering parameter inference challenging due to the inability to resolve individual parameter values along degenerate directions. Without knowledge of the degeneracy, we would not be able to predict how a change in one parameter would result in the same observations.

### 2.1 Pinpointing Candidate Parameters

**Mutual Information** To find degeneracies, we need to narrow down the parameter combinations that could cause combinatorial explosions as  $\binom{M}{k}$ . DEGENDETECTOR uses mutual information (MI) to rank tuples based on their statistical dependence. Recall that MI of a parameter pair  $I(\theta_a; \theta_b) = \int p(\theta_a, \theta_b) \log \frac{p(\theta_a, \theta_b)}{p(\theta_a)p(\theta_b)} d\theta_a d\theta_b$  measures the reduction in uncertainty about  $\theta_a$  given knowledge of  $\theta_b$ . If  $\theta_a$  and  $\theta_b$  are statistically independent, MI is zero; otherwise, it is strictly positive.

To estimate the pairwise MI scores, we use the  $k$ -nearest-neighbor estimator [Kraskov et al., 2004] and construct an  $M \times M$  matrix that is subsequently symmetrized and floored at zero. MI is sensitive to arbitrary nonlinear dependence structures between parameters, making it well suited to capture the nonlinear, multi-valued relationship responsible for physically motivated degeneracies.

For a target coupling depth of  $k$ , we enumerate all parameter tuples  $(\theta_{j_1}, \dots, \theta_{j_k})$ , assign them a score equal to the sum of all pairwise MI values among their elements, and arrange them in descending order. Higher aggregated MI implies a stronger, more significant degeneracy.

### 2.2 Symbolic Surface Fitting

We model the degeneracy as the level set of a separable function,

$$g_1(\theta_{j_1}) + g_2(\theta_{j_2}) + \dots + g_k(\theta_{j_k}) = c, \quad (2)$$

where  $g_\ell$ 's are univariate functions to be discovered and  $c$  is a real constant.

**Alternating Optimization** DEGENDETECTOR fits the component functions  $\{g_\ell\}_{\ell=1}^k$  and the constant  $c$  using alternating optimization, cycling through one component at a time while holding the others constant. This allows us to separate the  $k$ -dimensional fitting problem into  $k$  independent one-dimensional symbolic regression problems, resulting in a faster, more effective search. In what follows, we will use  $\theta_\ell$  as a shorthand for  $\theta_{j_\ell}$  for a fixed candidate tuple  $(j_1, \dots, j_k)$ .

To ensure that parameters with different units or physical scales contribute equally, we first normalize each parameter using the z-score  $\tilde{\theta}_\ell = \frac{\theta_\ell - \mu_\ell}{\sigma_\ell}$ , where  $\mu_\ell$  and  $\sigma_\ell$  are the sample mean and standard deviation. Each component function is then initialized to the identity,  $g_\ell(\tilde{\theta}_\ell) = \tilde{\theta}_\ell$ , and the constant to  $c = \text{mean}(\sum_\ell \tilde{\theta}_\ell)$ . At each alternating step, we cycle through all  $k$  components and establish the regression target:  $y_\ell = c - \sum_{m \neq \ell} g_m(\tilde{\theta}_m)$ , for component  $\ell$ , which is what  $g_\ell(\tilde{\theta}_\ell)$  ought to produce for every sample if the constraint is satisfied. We then perform symbolic regression on the pairs  $\{(g_\ell(\tilde{\theta}_\ell)^{(i)}, y_\ell^{(i)})\}_{i=1}^N$  using PYSR [Cranmer, 2023]. Once convergence is achieved, we use SYMPY [Meurer et al., 2017] to simplify and restore each component to the original scale, yielding a human-readable equation for the degeneracy.

**Orthogonal Fit Quality** For an implicit surface  $F(\boldsymbol{\theta}) = \sum_\ell g_\ell(\theta_\ell) - c = 0$ , the ordinary coefficient of determination  $R^2$  is ill-defined due to its asymmetry with respect to component function permutations (i.e., one parameter is specifically chosen to serve as the dependent variable). We define the

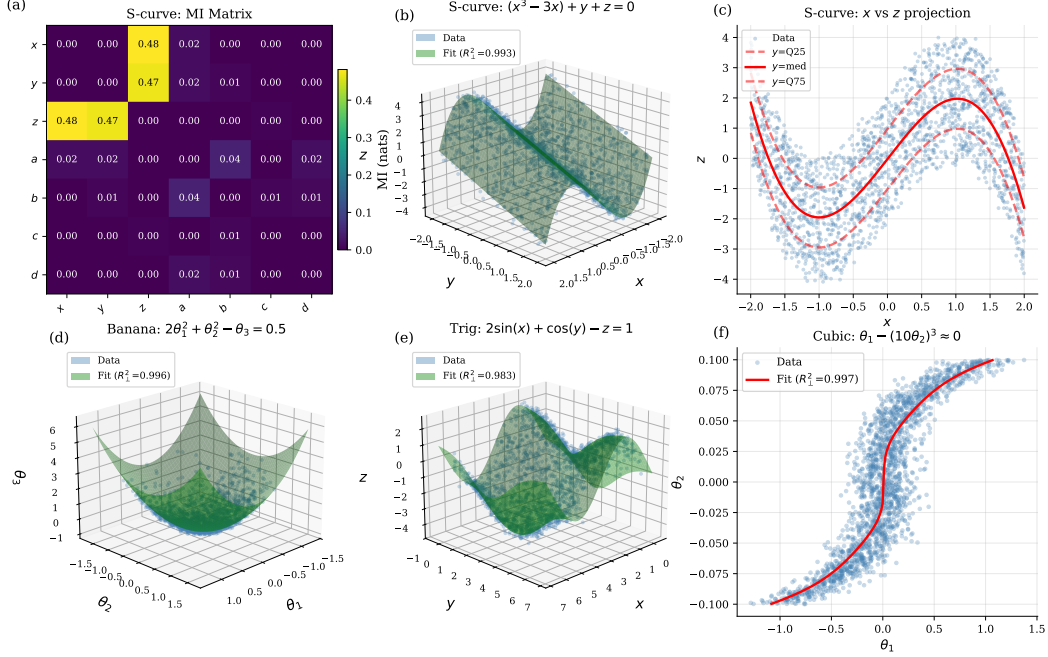


Figure 1: **Validation of DEGENDETECTOR on four benchmarks.** *Top row* shows the full pipeline for the S-curve degeneracy. (a) Pairwise MI matrix: the degenerate parameter tuple  $(x, y, z)$  forms a high MI cluster with pairwise MI  $\approx 0.47$ – $0.48$ , while nuisance parameters are cleanly separated. (b) Fitted surface (green) recovers the posterior samples with  $R^2_{\perp} = 0.9931$ . (c) Two-dimensional  $x$ - $z$  projection of the same fit. (d-f) shows fitted manifolds for the remaining three experiments.

orthogonal  $R^2$  as a more appropriate goodness of fit metric  $R^2_{\perp} = 1 - \mathcal{L}_{\perp}$  with the orthogonal loss being the mean squared perpendicular distance from the posterior samples to the fitted surface,

$$\mathcal{L}_{\perp} = \frac{1}{N} \sum_{i=1}^N \frac{[F(\boldsymbol{\theta}^{(i)})]^2}{\|\nabla F(\boldsymbol{\theta}^{(i)})\|^2} = \frac{1}{N} \sum_{i=1}^N \frac{[\sum_{\ell} g_{\ell}(\theta_{\ell}^{(i)}) - c]^2}{\sum_{\ell} [g'_{\ell}(\theta_{\ell}^{(i)})]^2 + \varepsilon}. \quad (3)$$

When true degeneracy is not separable, the orthogonal  $R^2$  will return a low score, indicating that the functional form has not been captured. The scientist can discover multiplicative degeneracies in log-space by activating LOGDEGEN, which applies a coordinate transformation  $\tilde{\theta}_i = \log \theta_i$ , reruns the pipeline on the transformed coordinates, and reports back in the original parameterization.

**Diagnostics** For each fitted tuple, DEGENDETECTOR provides a ranking of the top 5 candidate symbolic equations together with their  $R^2_{\perp}$  scores and expression complexities. A suite of diagnostic visualizations is also generated: a heatmap of the pairwise MI matrix identifying the most coupled parameter pairs, a corner plot of the degenerate parameters with the fitted surface overlaid as a contour, true-versus-predicted plots for each component function, and two- or three-dimensional renderings of the fitted manifold. These diagnostics allow the scientist to assess whether the separable model captures the true degeneracy structure before accepting the symbolic equation.

### 3 Benchmark Examples

We validate DEGENDETECTOR on synthetic posteriors that contain a priori defined, strong, non-Gaussian degeneracies along with some independent nuisance parameters (drawn from an isotropic Gaussian distribution). For the polynomial regime, we begin with the Rosenbrock function, a standard benchmark problem in numerical optimization [Goodman and Weare, 2010b] and extend the characteristic banana-shape constraints onto three dimensions. We then proceed to a more complex, non-monotone degeneracy in *S-curve*. In *Trig*, we assess if our method can approximate periodic

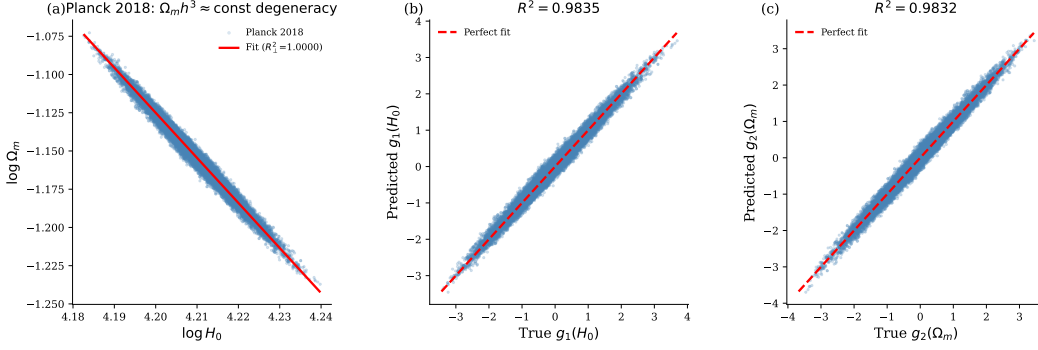


Figure 2: **Automated recovery of the CMB horizon-angle degeneracy from Planck 2018 posterior samples.** (a) Scatter of 25,225 Planck TTTEEE+lowl+lowE+lensing samples in log-parameter space, overlaid with the fitted degeneracy line (red). The tight linear locus confirms that  $\Omega_m h^3 \approx \text{const}$  is the dominant constraint. (b–c) True vs. predicted values for each component function  $g_1(H_0)$  and  $g_2(\Omega_m)$ , demonstrating that both are well described by simple log-linear forms.

functions with a limited set of polynomial-style operators, possibly through Taylor expansions or Gaussian-envelope estimations, and still achieve sufficient accuracy for the diagnostic to be useful. Finally, in *Cubic*, we test on posteriors sampled from an anisotropic prior with mismatched parameter scales, which ubiquitously occurs when physical parameters have different units.

Figure 1 illustrates how DEGENDETECTOR successfully recovers these non-trivial degeneracy across the four experiments. We find that DEGENDETECTOR can distinguish the underlying degenerate parameters from the other  $\binom{7}{3} = 35$  candidate tuples and recover an acceptable functional form (with  $R_{\perp}^2 > 0.98$ ).

## 4 Science Experiment

We apply DEGENDETECTOR to posterior chains from the Planck 2018 baseline analysis [Aghanim et al., 2020] with seven cosmological parameters  $\Omega_b h^2$ ,  $\Omega_c h^2$ ,  $H_0$ ,  $\Omega_m$ ,  $\sigma_8$ ,  $n_s$ , and  $\tau$ . The CMB horizon-angle degeneracy  $\Omega_m h^3 \approx \text{const}$  is well established; the exponent of 3 [Percival et al., 2002, empirically find  $\sim 3.4$ ] reflects the combined dependence of the angular scale of the sound horizon at last scattering on the matter density and expansion rate [Efstathiou and Bond, 1999]. Although the power law degeneracy is non-separable in the  $(\Omega_m, h)$  space, it can be recovered as the level set  $\log \Omega_m + 3 \log h = \text{const}$  in  $(\log \Omega_m, \log H_0)$  space.

Without any prior knowledge of the physics, LOGDEGEN recovers  $(H_0, \Omega_m)$ , with a mutual-information score of  $I = 2.07$  and fits the Planck posterior in log-space, returning the expression:

$$123.97 \log H_0 + 42.07 \log \Omega_m = \text{const}, \quad (4)$$

with  $R_{\perp}^2 \approx 0.98$  and a residual standard deviation of  $\sigma_{\perp} = 0.128$  in log-space. Dividing by the  $\Omega_m$  coefficient yields a ratio of  $123.97/42.07 \approx 2.947$ , which is reasonably consistent with the anticipated functional form up to a constant. The manifold fit and its agreement with the posterior samples are illustrated in Figure 2.

## 5 Conclusion

We developed DEGENDETECTOR, a novel degeneracy detection framework that integrates mutual information ranking with alternating symbolic regression to identify coupled parameters and render that degeneracy as a comprehensible equation. On synthetic benchmarks and Planck 2018 posteriors, we recover the true functional form of degeneracies with  $R_{\perp}^2 > 0.98$  across all cases.

The primary limitation in our method is the assumption of separability; future research could extend the framework to non-separable functional forms via multivariate symbolic regression or learned reparameterizations.

## References

- N. Aghanim et al. Planck 2018 results. VI. Cosmological parameters. , 641:A6, 2020. doi: 10.1051/0004-6361/201833910.
- Eoin Ó Colgáin, Saeed Pourjaghi, M. M. Sheikh-Jabbari, and Darragh Sherwin. A comparison of bayesian and frequentist confidence intervals in the presence of a late universe degeneracy, 2025. URL <https://arxiv.org/abs/2307.16349>.
- Miles Cranmer. Interpretable machine learning for science with PySR and SymbolicRegression.jl, 2023.
- George Efstathiou and J. R. Bond. Cosmic confusion: degeneracies among cosmological parameters derived from measurements of microwave background anisotropies. , 304(1):75–97, 1999. doi: 10.1046/j.1365-8711.1999.02274.x.
- Janis Fluri, Tomasz Kacprzak, Alexandre Refregier, Aurelien Lucchi, and Thomas Hofmann. Cosmological parameter estimation and inference using deep summaries. *Physical Review D*, 104(12), December 2021. ISSN 2470-0029. doi: 10.1103/physrevd.104.123526. URL <http://dx.doi.org/10.1103/PhysRevD.104.123526>.
- Jonathan Goodman and Jonathan Weare. Ensemble samplers with affine invariance. *Communications in Applied Mathematics and Computational Science*, 5(1):65–80, 2010a. doi: 10.2140/camcos.2010.5.65.
- Jonathan Goodman and Jonathan Weare. Ensemble samplers with affine invariance. *Communications in applied mathematics and computational science*, 5(1):65–80, 2010b.
- Bhuvnesh Jain and Uroš Seljak. Cosmological model predictions for weak lensing: Linear and nonlinear regimes. , 484(2):560–573, 1997. doi: 10.1086/304372.
- Jens Jasche and Benjamin D. Wandelt. Bayesian physical reconstruction of initial conditions from large-scale structure surveys. , 432(2):894–913, 2013. doi: 10.1093/mnras/stt449.
- Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical Review E*, 69(6):066138, 2004. doi: 10.1103/PhysRevE.69.066138.
- Aaron Meurer et al. SymPy: symbolic computing in Python. *PeerJ Computer Science*, 3:e103, 2017. doi: 10.7717/peerj-cs.103.
- George Papamakarios, David C. Sterratt, and Iain Murray. Sequential neural likelihood: Fast likelihood-free inference with autoregressive flows. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 837–848, 2019.
- W. J. Percival, W. Sutherland, J. A. Peacock, C. M. Baugh, J. Bland-Hawthorn, T. Bridges, R. Cannon, S. Cole, M. Colless, C. Collins, W. Couch, G. Dalton, R. De Propris, S. P. Driver, G. Efstathiou, R. S. Ellis, C. S. Frenk, K. Glazebrook, C. Jackson, O. Lahav, I. Lewis, S. Lumsden, S. Maddox, S. Moody, P. Norberg, B. A. Peterson, and K. Taylor. Parameter constraints for flat cosmologies from cosmic microwave background and 2dfgrs power spectra. *Monthly Notices of the Royal Astronomical Society*, 337(3):1068–1080, December 2002. ISSN 1365-2966. doi: 10.1046/j.1365-8711.2002.06001.x. URL <http://dx.doi.org/10.1046/j.1365-8711.2002.06001.x>.
- John Skilling. Nested sampling for general Bayesian computation. *Bayesian Analysis*, 1(4):833–859, 2006. doi: 10.1214/06-BA127.
- Eleonora Villa, Luigi D’Amico, Aldo Barca, Fatima Modica Bittordo, Francesco Ali, Massimo Meneghetti, and Luca Naso. Addressing prior dependence in hierarchical bayesian modeling for pta data analysis ii: Noise and sgwb inference through parameter decorrelation, 2025. URL <https://arxiv.org/abs/2511.01959>.