

Hypothesis Tests Under Separation*

Carlisle Rainey[†]

Version from June 3, 2022.

Separation commonly occurs in political science, usually when a binary explanatory variable perfectly predicts a binary outcome. In these situations, methodologists often recommend penalized maximum likelihood or Bayesian estimation. But researchers might struggle to identify an appropriate penalty or prior distribution. I show that while the popular Wald test produces misleading (even nonsensical) p -values under separation, likelihood ratio tests and score tests behave in the usual manner. Therefore, researchers can produce meaningful p -values with standard frequentist tools under separation *without the use of penalties or prior information*.

*All data and code for the paper are available at <https://github.com/carlislerainey/wilks/>. A computational companion that illustrates how one can compute the quantities I discuss in the paper is available at <https://github.com/carlislerainey/wilks/blob/master/doc/cc-wilks.md>.

[†]Carlisle Rainey is Associate Professor of Political Science, Florida State University, 540 Bellamy, Tallahassee, FL, 32306. (crainey@fsu.edu).

Separation commonly occurs in political science, usually when a binary explanatory variable perfectly predicts a binary outcome (e.g., Gustafson 2020; Mehltritt 2021; Owsiak and Vasquez 2021). For example, Barrilleaux and Rainey (2014) find that being a Democrat perfectly predicts a governor supporting Medicaid expansion under the Affordable Care Act. Under separation, the usual maximum likelihood estimates are unreasonably large and the Wald p -values are highly misleading.

As a solution, some methodologists propose using a Bayesian prior distribution to regularize the estimates, which we can alternatively consider as a *penalized* maximum likelihood estimator. Zorn (2005; see also Heinze and Schemper 2002) points political scientists toward the penalized maximum likelihood estimator (Firth 1993), which is equivalent to Jeffreys prior distribution (Jeffreys 1946), while Gelman et al. (2008) recommend a Cauchy prior distribution. Both methods ensure finite estimates in theory and usually produce reasonably-sized estimates in practice. Methodologists continue to recommend these penalized or Bayesian estimators as a solution to separation (e.g., Cook, Niehaus, and Zuhlke 2018; Anderson, Bagozzi, and Koren 2021; Cook, Hays, and Franzese 2020; Crisman-Cox, Gasparyan, and Signorino 2022).

But Rainey (2016) points out that the estimates (and especially the confidence intervals) depend largely on the chosen prior or penalty. Many priors guarantee finite estimates and lead to meaningfully different conclusions. He argues that the set of *a priori* “reasonable” and “implausible” parameters depends on the substantive application, so context-free defaults (like Jeffreys and Cauchy priors) might not produce reasonable results. Starkly emphasizing this point, Beiser-McGrath (2022) shows that Jeffreys prior can lead to (statistically significant) estimates in the *opposite direction* of the separation. Rainey (2016) concludes that “[w]hen facing separation, researchers must *carefully* choose a prior distribution to nearly rule out implausibly large effects” (p. 354). But it is not always easy to include prior information, and some scholars prefer to avoid injecting prior information into their model. How can researchers proceed in these situations? In particular, can

they obtain useful p -values to test hypotheses in the usual frequentist framework without using prior information?

I show that while the popular Wald test produces misleading (even nonsensical) p -values under separation, likelihood ratio tests and score tests behave in the usual manner. As such, researchers can produce meaningful p -values with standard frequentist tools under separation *without the use of prior information*.

Hypothesis Tests Under Separation

Maximum likelihood provides a general and powerful framework for obtaining estimates of model parameters and testing hypotheses. In our case of logistic regression, we write the probability π_i that an event occurs for observation i of n (or that the outcome variable $y_i = 1$) as

$$\pi_i = \text{logit}^{-1}(X_i\beta) \text{ for } i = 1, 2, \dots, n,$$

where X represents a matrix of explanatory variables and β represents a vector of coefficients. Then we have the likelihood function

$$L(\beta|y) = \pi_i^{y_i}(1 - \pi_i)^{(1-y_i)}$$

and the log-likelihood function

$$\ell(\beta|y) = \log L(\beta|y) = y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i).$$

Researchers typically use numerical algorithms to locate the maximum likelihood estimate $\hat{\beta}^{ML}$ that maximizes ℓ . Then the researchers use certain features of ℓ to test hypotheses. To fix ideas, I focus on the point null hypothesis $H_0 : \beta_s = 0$. However, the intuition and conclusions generalize to more complex hypotheses.

The literature offers three common methods to assess the null hypothesis—the “holy trinity” of hypothesis tests: the Wald test, the likelihood ratio test, and the score test. For practical reasons, most regression tables in political science report Wald p -values. However, the Wald test is uniquely ill-suited for testing hypotheses under separation. Because the Wald test fails, some researchers (too) quickly turn to penalized and Bayesian estimators. However, the usual likelihood ratio and score tests work as expected under separation. Below I briefly describe each test, explain why the Wald test works poorly under separation, and describe why the likelihood ratio and score tests perform better.

Wald Test

The Wald test uses the shape of the log-likelihood function around the maximum to estimate the precision of the point estimate. If small changes in the parameter near the maximum lead to large changes in the log-likelihood function, then we can treat the maximum likelihood estimate as precise. We usually estimate the standard error $\widehat{SE}(\hat{\beta}_i^{ML})$ as

$$\widehat{SE}(\hat{\beta}_i^{ML}) = \left(-\frac{\partial^2 \ell(\hat{\beta}_i^{ML}|y)}{\partial^2 \hat{\beta}_i^{ML}} \right)^{-\frac{1}{2}}.$$

Wald (1943) advises us how compare the estimate with the standard error: the statistic $Z_w = \frac{\hat{\beta}_i^{ML}}{\widehat{SE}(\hat{\beta}_i^{ML})}$ approximately follows a standard normal distribution (Casella and Berger 2003, pp. 492-493; Greene 2012, pp. 527-529).

The Wald approach works poorly when dealing with separation. Under separation, the log-likelihood function at the numerical maximum is nearly flat. The flatness produces very large standard error estimates—much larger than the coefficient estimates. Figure 1 shows this intuition for a typical, non-monotonic log-likelihood function (i.e., without separation) and a monotonic log-likelihood function (i.e., with separation). In the absence of separation, the curvature of the log-likelihood function around the maximum speaks to

the evidence against the null hypothesis. But under separation, the monotonic likelihood function is flat at the maximum, regardless of the relative likelihood of the data under the null hypothesis.

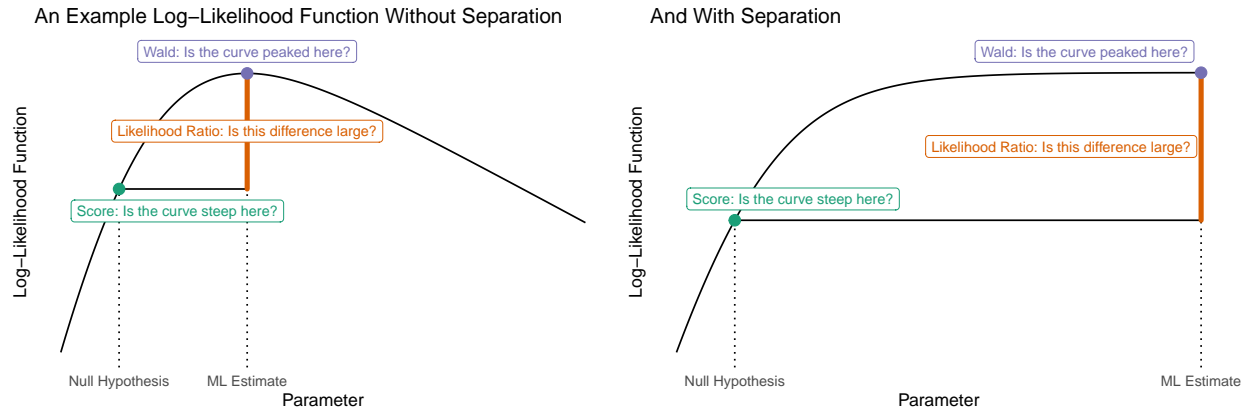


Figure 1: A figure summarizing logic of the "holy trinity" of hypothesis tests. The Wald test relies on the curvature around the maximum of the log-likelihood function, which breaks down under separation. But the likelihood ratio and score test rely on *other* features of the log-likelihood function that are not meaningfully impacted by separation.

We can develop the intuition of the more precisely and formally. Suppose that a binary explanatory variable s with coefficient β_s perfectly predicts the outcome y_i such that when $s_i = 1$ then $y_i = 1$. Then the log-likelihood function increases in β_s . The standard error estimate associated with each β_s increases as well. But critically, the estimated standard error increases *faster* than the associated coefficient, because $\lim_{\beta_s \rightarrow \infty} \left[\left(-\frac{\partial^2 \ell(\beta_s | y)}{\partial^2 \beta_s} \right)^{-\frac{1}{2}} - \beta_s \right] = \infty$. Thus, under separation, the estimated standard error will be much larger than the coefficient for the separating variable for any algorithm that obtains a sufficiently large coefficient. This implies two conclusions. First, so long as the researcher uses a sufficiently precise algorithm, *the Wald test will never reject the null hypothesis*, regardless of the data set. Second, if the Wald test can never reject the null hypothesis for any data set with separation, then the power of the test is strictly bounded by the chance of separation. In particular, *the power of the test cannot exceed* $1 - \Pr(\text{separation})$. If the data set features separation in nearly 100% of repeated samples,

then the Wald test will have power near 0%.

As a final illustration, suppose an absurd example in which a binary treatment perfectly predicts 500 successes and 500 failures (i.e., $y = x$ always). Of course, this data set is *extremely* unlikely under the null hypothesis that the coefficient for the treatment indicator equals zero. The exact p -value for the null hypothesis that successes and failures are equally likely under both treatment and control equals $2 \times \left(\frac{1}{2}\right)^{500} \times \left(\frac{1}{2}\right)^{500} = \frac{2}{2^{1000}} \approx \frac{2}{10^{301}}$. (For comparison, there are about 10^{80} atoms in the known universe.) Yet, the default `glm()` routine in R calculates a Wald p -value of 0.998 with the default precision (and 1.000 with the maximum precision). When dealing with separation, the Wald test breaks down; researchers cannot use the Wald test to obtain reasonable p -values for the coefficient of a separating variable.

Likelihood Ratio Test

The likelihood ratio test resolves the problem of the flat log-likelihood by comparing the maximum log-likelihood of two models: an “unrestricted” model ML that imposes no bounds on the estimates and a “restricted” model ML_0 that constrains the estimates to the region suggested by the null hypothesis. If the data set is much more likely under the unrestricted estimate than under the restricted estimate, then the researcher can reject the null hypothesis. Wilks (1938) advises us how to compare the unrestricted log-likelihood $\ell(\hat{\beta}^{ML}|y)$ to the restricted log-likelihood $\ell(\hat{\beta}^{ML_0}|y)$: $D = 2 \times [\ell(\hat{\beta}^{ML}|y) - \ell(\hat{\beta}^{ML_0}|y)]$ approximately follows a χ^2 distribution with degrees of freedom equal to the number of constrained dimensions (Casella and Berger 2003, pp. 488-492, Greene 2012, pp. 526-527).

Figure 1 shows the intuition of the likelihood ratio test. The gap between the unrestricted and restricted maximum summarizes the evidence against the null hypothesis. Importantly, the logic does not break down under separation. Unlike the Wald test, the likelihood ratio test *can* reject the null hypothesis under separation.

Score Test

The score test resolves the problem of the flat log-likelihood by evaluating the gradient of the log-likelihood function at the null hypothesis. If the log-likelihood function is increasing rapidly at the null hypothesis, this casts doubt on the null hypothesis. The score test uses the score function $S(\beta) = \frac{\partial \ell(\beta|y)}{\partial \beta}$ and the Fisher information $I(\beta) = -E_{\beta} \left(\frac{\partial^2 \ell(\beta|y)}{\partial^2 \beta} \right)$. When evaluated at the null hypothesis, the score function quantifies the slope and the Fisher information quantifies the variance of that slope in repeated samples. If the score at the null hypothesis is large, then the researcher can reject the null hypothesis. Rao (1948) advises us how to compare the score to its standard error: $Z_s = \frac{S(\beta_s^0)}{\sqrt{I(\beta_s^0)}}$ follows a standard normal distribution (Casella and Berger 2003, pp. 494-495, Greene 2012, pp. 529-530).

Figure 1 shows the intuition of the score test. The slope of the log-likelihood function under the null hypothesis summarizes the evidence against the null hypothesis. As with the likelihood ratio test, the logic works even under separation, and the score test *can* reject the null hypothesis under separation.

Table 1 summarizes the three tests. Most importantly, the likelihood ratio and score tests rely on features of the log-likelihood function that are not meaningfully affected by a monotonic log-likelihood function. The Wald test, on the other hand, cannot provide a reasonable test under separation.

Table 1: A table summarizing the "holy trinity" of hypothesis tests.

Test	Feature	Statistic and Distribution
Wald	Curvature of the log-likelihood function around the maximum.	$Z_w = \frac{\hat{\beta}_i^{ML}}{\widehat{SE}(\hat{\beta}_i^{ML})}$ follows a standard normal distribution.
Likelihood Ratio	Relative log-likelihoods of the unrestricted and restricted models.	$D = 2 \times [\ell(\hat{\beta}^{ML} y) - \ell(\hat{\beta}^{ML_0} y)]$ follows a χ^2 distribution with degrees of freedom equal to the number of constrained dimensions.
Score	Slope of the log-likelihood function <i>at the null hypothesis</i> .	$Z_s = \frac{S(\beta_s^0)}{\sqrt{I(\beta_s^0)}}$ follows a standard normal distribution.

Simulations

To evaluate the performance of the three tests under separation, I use a diverse collection of data-generating processes (DGPs). Each DGP features separation in more than 10% of repeated samples. Importantly, I cannot focus on data sets *with separation* because separation is a feature of a particular sample. Instead, I focus on DGPs that *sometimes* feature separation (e.g., in 15% of repeated samples, in 50% of repeated samples, etc.).

To create the collection of DGPs, I imagine the logistic regression model $\Pr(y = 1) = \text{logit}^{-1}(\beta_{\text{cons}} + \beta_s s + \beta_{z_1} z_1 + \dots + \beta_{z_k} z_k)$ and a researcher testing the null hypothesis that the binary explanatory variable s (that might produce separation) has no effect on a binary outcome variable y (i.e., that $\beta_s = 0$). I vary the total number of observations, the frequency that $s = 1$, the value of β_{cons} , and the number of control variables (k).¹ For each DGP, I use Monte Carlo simulation to compute the power function for each of the three tests.

A Close Look at a Single DGP

Following the structure discussed above, I first examine a single DGP. There are 50 total observations, $s = 1$ for five of the observations and $s = 0$ for the other 45 observations, the intercept β_{cons} is zero, and there are two control variables. Table 2 shows the power function for each of the three tests and the percent of the data sets that featured separation for the parameter combination and the value of β_s . For this particular DGP, separation is relatively rare when β_s —the coefficient for the potentially separating variable—is near zero. But for $\beta_s = \pm 2$, about 50% of the data sets feature separation. And for β_s larger/smaller than ± 4 , more than 90% of the data sets feature separation.

This DGP clearly demonstrates the poor performance of the Wald test. Even though the data sets with separation should allow the researcher to reject the null hypothesis, at

¹Each control variable z_i is simulated from a normal distribution with a standard deviation of 0.5. The coefficient for each control variable β_{z_i} is set to one.

Table 2

β_s	Ideal Power	Percent with Separation	Wald Test Power	Likelihood Ratio Test Power	Score Test Power
-10	As high as possible.	100%	0%	99%	85%
-7		99%	0%	99%	85%
-4		89%	0%	90%	77%
-2		45%	2%	49%	42%
-1		13%	2%	18%	15%
0	5%	4%	1%	6%	5%
1	As high as possible.	19%	0%	20%	16%
2		50%	0%	51%	41%
4		90%	0%	89%	70%
7		99%	0%	97%	76%
10		100%	0%	98%	78%

This table shows the power for the Wald, likelihood ratio, and score tests for a DGP that often features separation.

least occasionally, the power of the Wald test is near 0% even for very large effects. This happens because the Wald test cannot reject the null hypothesis under separation. When separation happen often, the power *must* be low. The likelihood ratio and score tests, on the other hand, perform as expected. For both alternatives, the power of the test when $\beta_s = 0$ is about 5%, as designed, and the power approaches 100% relatively quickly as β_s moves away from zero.

A Broad Look at Many DGPs

Table 3 shows the broad collection of parameters I used in the simulations. All combinations of the values below yield 150 unique combinations. I exclude combinations in which the frequency that $s = 1$ exceeded the number of observations, which leaves 110 combinations. For each of these 110 combinations, I use Monte Carlo simulations to estimate the probability that each method rejects the null hypothesis as the coefficient β_s varies from -10 to 10 (across the particular values shown in Table 2). I also exclude any scenario in which the parameter and the coefficient β_s would sometimes produce data sets with no variation in the outcome variable. As such, one might consider this a diverse collection of DGPs that sometimes generate separation but almost always have variation in the outcome.

Figure 2 shows the statistical power of each of the three tests as the chance of

Table 3

Parameter	Value
Total Number of Observations	50, 250, 500
Frequency that $s = 1$	5, 25, 50, 100, 250
The Value of the Intercept β_{cons}	-8, -5, -2.5, -1, 0
The Number of Control Variables k	2, 6

separation varies across the many DGPs. Most starkly, the power of the Wald test is bounded above by $1 - \Pr(\text{separation})$, and many scenarios achieve the boundary. Intuitively, as the chance of separation increases, the power of the test should increase as well, because separation is evidence of a large coefficient. But because a large coefficient makes separation more likely, a large coefficient *decreases* the power of the Wald test. The likelihood ratio and score test, on the other hand, behave as expected.

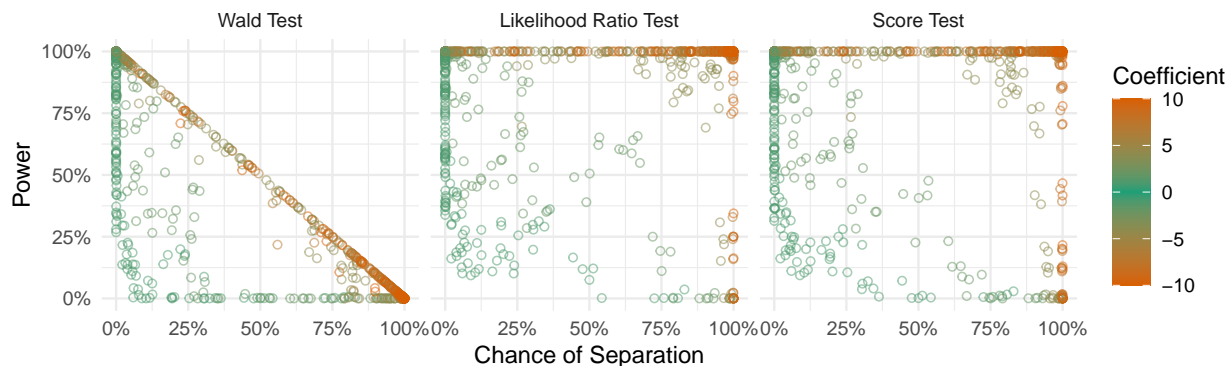


Figure 2: The power of tests across a range of scenarios as the chance of separation varies.

Figure 3 plots the power function for each of the 110 DGPs in the diverse collection.² The fine lines show the power function for each DGP and the color of the points indicates the chance of separation at that point. The heavy, solid lines show the average of the collection of power functions. The heavy, dashed lines show the power functions for the other two tests, for comparison.

First, the power functions for the Wald tests show its poor properties. For most of the

²I exclude parameter combinations that might not produce variation in the outcome variable, so I do not compute the power function across the entire range of the coefficient for every power function. For 47 of the 110 DGPs, I exclude a portion of the range. For 63 of the 110, I compute the power function across the entire range from -10 to 10.

power functions, as the true coefficient grows larger in magnitude from about three, the test becomes less powerful. This occurs because separation becomes more likely and the test cannot reject the null hypothesis when separation occurs. Second, the likelihood ratio and score test behave reasonably well. Most importantly, the size of the likelihood ratio and score tests is about 5% when the coefficient equals zero and grows as the coefficient moves away from zero. The likelihood ratio test is slightly preferred, at least of this collection of DGPs, because it is slightly more powerful when the coefficient does not equal zero, especially for large, negative coefficients.³

Re-Analysis of Barrilleaux and Rainey (2014)

To illustrate the power and simplicity of frequentist hypothesis testing under separation, I reanalyze data from Barrilleaux and Rainey (2014), who examine U.S. state governors decisions to support or oppose the Medicaid expansion under the 2010 Affordable Care Act. But because all Democratic governors supported the expansion, separation occurs—being a Democratic governor perfectly predicts support for Medicaid expansion.

I focus on their first hypothesis: *Republican governors are more likely to oppose the Medicaid expansion funds than Democratic governors*. Barrilleaux and Rainey adopt a fully Bayesian approach, modeling the probability that a state's governor opposes the Medicaid expansion as a function of the governor's partisanship and several other covariates. Here, I re-estimate their logistic regression model and compute Wald, likelihood ratio, and score p -values. Table 4 presents the estimates and p -values for the coefficient of the binary Democratic governor variable. While the Wald p -value is near 1.00, the likelihood ratio and score tests produce more plausible p -values of 0.003 and 0.009, respectively.⁴

³This asymmetry occurs because I only consider negative intercepts.

⁴For comparison, the Wald tests for the Jeffreys and Cauchy penalized maximum likelihood estimators produce p -values of 0.060 and 0.038, respectively.

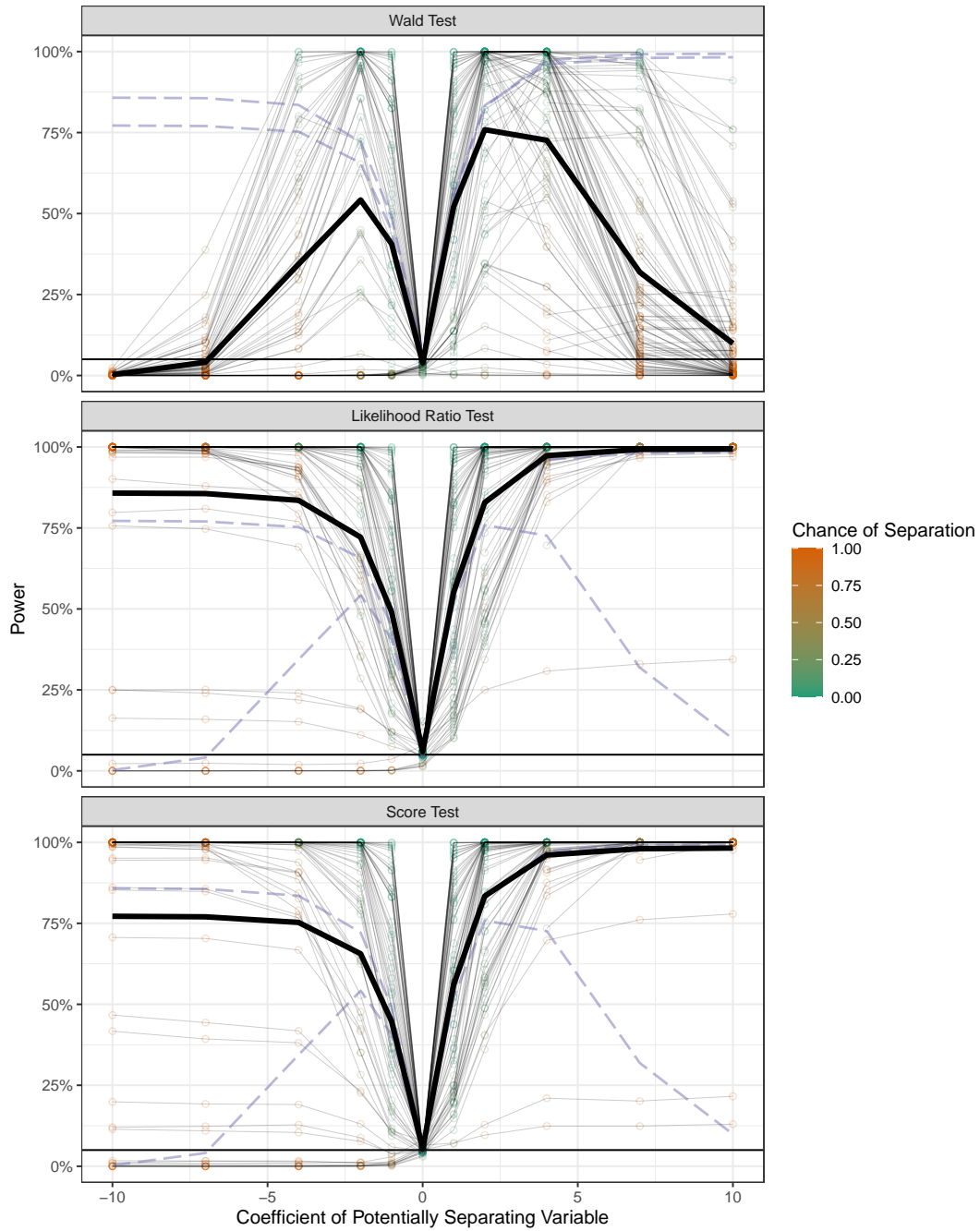


Figure 3: The power functions for the Wald, likelihood ratio, and score tests for a diverse collection of data-generating processes that sometimes generate separation.

Conclusion

Separation commonly occurs in political science. When this happens, I show that the usual Wald p -values are highly misleading. But researchers cannot always use prior

Table 4

Estimator	Coef. Est.	SE Est.	Wald <i>p</i> -Value	LR <i>p</i> -Value	Score <i>p</i> -Value
ML with Default Precision	-20.349	3,224	0.995	0.003	0.009
ML with Maximum Precision	-35.223	15 million	1.000	0.003	0.009

The *p*-values from several procedures that researchers might use when dealing with separation in logistic regression models. The Wald test for maximum likelihood estimates relies on unreasonable standard errors that depend heavily on the precision of the algorithm and, as a consequence, produces unrealistic *p*-values. However, the likelihood ratio and score tests produce reasonable *p*-values.

information to address a monotonic likelihood function. Even without a suitable prior or penalty, I show that the standard likelihood ratio and score tests behave in the usual way. As such, researchers can use the likelihood ratio and score tests to produce meaningful *p*-values under separation *without the use of prior information*.

References

- Anderson, Noel, Benjamin E Bagozzi, and Ore Koren. 2021. "Addressing Monotone Likelihood in Duration Modelling of Political Events." *British Journal of Political Science* 51(4): 1654–71.
- Barrilleaux, Charles, and Carlisle Rainey. 2014. "The Politics of Need: Examining Governors' Decisions to Oppose the 'Obamacare' Medicaid Expansion." *State Politics and Policy Quarterly* 14(4): 437–60.
- Beiser-McGrath, Liam F. 2022. "Separation and Rare Events." *Political Science Research and Methods* 10(2): 428–37.
- Casella, George, and Roger L. Berger. 2002. *Statistical Inference*. 2nd ed. Pacific Grove, CA: Duxbury.
- Cook, Scott J, Jude C Hays, and Robert J Franzese. 2020. "Fixed Effects in Rare Events Data: A Penalized Maximum Likelihood Solution." *Political Science Research and Methods* 8(1): 92–105.
- Cook, Scott J, John Niehaus, and Samantha Zuhlke. 2018. "A Warning on Separation in Multinomial Logistic Models." *Research & Politics* 5(2): 2053168018769510.
- Crisman-Cox, Casey, Olga Gasparyan, and Curtis S Signorino. 2022. "Finding and Accounting for Separation Bias in Strategic Choice Models."
- Firth, David. 1993. "Bias Reduction of Maximum Likelihood Estimates." *Biometrika* 80(1): 27–38.
- Gelman, Andrew, Aleks Jakulin, Maria Grazia Pittau, and Yu-Sung Su. 2008. "A Weakly Informative Prior Distribution for Logistic and Other Regression Models." *The Annals of Applied Statistics* 2(4): 1360–83.

- Greene, William H. 2012. *Econometric Analysis*. 7th ed. Upper Saddle River, New Jersey: Prentice Hall.
- Gustafson, Daniel. 2020. "Hunger to Violence: Explaining the Violent Escalation of Nonviolent Demonstrations." *Journal of Conflict Resolution* 64(6): 1121–45.
- Heinze, Georg, and Michael Schemper. 2002. "A Solution to the Problem of Separation in Logistic Regression." *Statistics in Medicine* 21(16): 2409–19.
- Jeffreys, H. 1946. "An Invariant Form of the Prior Probability in Estimation Problems." *Proceedings of the Royal Society of London, Series A* 186(1007): 453–61.
- Mehrtretter, Andreas. 2021. "Arming for Conflict, Arming for Peace? How Small Arms Imports Affect Intrastate Conflict Risk." *Conflict Management and Peace Science*: Forthcoming.
- Owsiak, Andrew P, and John A Vasquez. 2021. "Peaceful Dyads: A Territorial Perspective." *International Interactions* 47(6): 1040–68.
- Rainey, Carlisle. 2016. "Dealing with Separation in Logistic Regression Models." *Political Analysis* 24(3): 339–55.
- Rao, C Radhakrishna. 1948. "Large Sample Tests of Statistical Hypotheses Concerning Several Parameters with Applications to Problems of Estimation." In *Mathematical Proceedings of the Cambridge Philosophical Society*, Cambridge University Press, 50–57.
- Wald, Abraham. 1943. "Tests of Statistical Hypotheses Concerning Several Parameters When the Number of Observations Is Large." *Transactions of the American Mathematical society* 54(3): 426–82.
- Wilks, Samuel S. 1938. "The Large-Sample Distribution of the Likelihood Ratio for Testing Composite Hypotheses." *The annals of mathematical statistics* 9(1): 60–62.
- Zorn, Christopher. 2005. "A Solution to Separation in Binary Response Models." *Political*

Analysis 13(2): 157–70.