

Inference in generalized linear models with robustness to misspecified variances

Riccardo De Santis^{*1}, Jelle J. Goeman², Jesse Hemerik³, Samuel Davenport⁴, and Livio Finos⁵

¹University of Siena, Italy

²Leiden University Medical Center, The Netherlands

³Erasmus University Rotterdam, The Netherlands

⁴University of California San Diego, United States

⁵University of Padova, Italy

Abstract

Generalized linear models usually assume a common dispersion parameter, an assumption that is seldom true in practice. Consequently, standard parametric methods may suffer appreciable loss of type I error control. As an alternative, we present a semi-parametric group-invariance method based on sign flipping of score contributions. Our method requires only the correct specification of the mean model, but is robust against any misspecification of the variance. We present tests for single as well as multiple regression coefficients. The test is asymptotically valid but shows excellent performance in small samples. We illustrate the method using RNA sequencing count data, for which it is difficult to model the overdispersion correctly. The method is available in the R library `flipscores`.

1 Introduction

When testing for equality of means between two groups, statisticians often worry that there may be a difference in variance between the two groups. The task of properly taking this variance into account, known as the Behrens-Fisher problem (Fisher, 1935, 1941), has generated a huge statistical literature (Kim and Cohen, 1998; Chang and Pal, 2008). In the linear regression model (LM), which generalizes the two-group model, the assumption that all error terms have equal variance is less often questioned. While goodness-of-fit tests and other diagnostics for heteroscedasticity exist (Goldfeld and Quandt, 1965; Glejser,

^{*}To contact: riccardo.desantis2@unisi.it

1969; Breusch and Pagan, 1979; Jarque and Bera, 1980; Cook and Weisberg, 1983; Long and Ervin, 2000), there is no simple and general way to follow up on obvious lack of fit (Rochon et al., 2012). Still, the problems arising from misspecified variances in regression can be as severe in regression as they are for the two-group model, especially when the variance depends on the covariates in the model in an unknown way: standard errors are too large or too small, and statistical tests can become severely conservative or anti-conservative.

The situation is worse if we broaden the perspective to generalized linear models (GLMs). Overdispersed GLMs, that allow for additional variation between subjects, have become important in many fields; e.g. negative binomial or quasi-Poisson models in RNA sequencing (Love et al., 2014). GLMs and overdispersed GLMs generally use a single dispersion parameter for all subjects. This assumption is similar to the common variance assumption in the linear model, and violation of the assumption can lead to loss of type I error control in the same way. For large data sets this problem can be addressed by estimating the variance of the test statistic robustly using the sandwich estimator (Huber, 1967; White, 1982) or by direct modelling using Generalized Additive Models for Location Scale and Shape (Rigby and Stasinopoulos, 2005). However, these methods can have poor error control for small sample sizes (Boos, 1992; Freedman, 2006; Maas and Hox, 2004; Kauermann and Carroll, 2000).

In linear models the wild bootstrap (Davidson and Flachaire, 2008) addresses the problem of unknown heteroscedasticity in a different way by randomly sign-flipping transformed residuals. This method bypasses the need to model the variance by comparing each subject’s contribution to the test statistic with its own sign-flipped copy only. The method of Hemerik et al. (2020) extends this method to generalized linear models and to situations with nuisance parameters by sign-flipping the contributions to the score statistic rather than the residuals. However, their contribution has two main drawbacks. First, the proven robustness of the methods to variance misspecification is rather limited, relating only to misspecification by a constant. Second, although it has better control of type I error than Sandwich-based methods, it still struggles to control type I error properly for small sample sizes.

In this paper, we present an improvement of the method of Hemerik et al. (2020), addressing both drawbacks. In the first place we develop a standardization of the sign-flip-based test statistics that boosts the convergence to the asymptotic distribution and greatly improves the control of the type one error. Our proposed standardization is innovative because the standardization factor is conditional on the random sign-flips and does not require model refitting as is common for resampling based standardized test statistics (e.g., methods reviewed in Winkler et al., 2014; Frossard and Renaud, 2021). Second, we prove robustness to any misspecification of the variance, greatly improving the scope of the method.

The structure of this paper is as follows. In Section 2 we summarize the main results of the paper, placing them in the context of the literature. In Section 3 we introduce the modeling context and assumptions that we will consider. After recalling the existing score flipping test in Section 4, in Section 5 we propose

the novel test. Section 6 contains the most important results of this paper, since here we prove the robustness properties of our tests. An extension to multivariate tests is in Section 7. Section 8 contains the simulation study and Section 9 a data analysis, which illustrates the relevance of our approach to the analysis of RNA-Seq data.

2 Contributions of the paper

We start by summarizing the two main novel contributions of this paper, placing them in the context of the literature.

2.1 Robustness to variance misspecification in GLMs

While many methods have been proposed to remedy the problem variance misspecification, they all come at a price.

A general and popular approach that is robust to variance misspecification is the sandwich estimator (Cox, 1961, 1962; Eicker, 1967; Huber, 1967; White, 1982). However, this method is well-known to have poor type I error control for small sample sizes (Boos, 1992; Kauermann and Carroll, 2000; Maas and Hox, 2004; Freedman, 2006). Alternatively, more flexible models can be used to tackle the problem of the strict assumptions made by the GLMs. Generalized additive models for location, scale and shape (GAMLSS; Rigby and Stasinopoulos, 2005) permit to jointly model the mean and the variance as a function of covariates. However, these models tend to be overparametrized, which can also result in poor small sample performance, as we will see in the simulations of Section 8. Furthermore, constructing and evaluating a variance model (GAMLSS; Rigby and Stasinopoulos, 2005) is more complex in GLMs than in the linear model, since the variance of the residuals differs between observations even in a correctly specified model.

A completely different perspective is taken by HulC (Hull based Confidence) (Kuchibhotla et al., 2023), a recent robust semiparametric method based on sample splitting, whose key assumption is only the asymptotic median unbiasedness of the estimator. This method gives confidence intervals with theoretical guarantees. However, as a sample splitting method, it requires a minimal sample size, and comes with a relevant loss of power. We will consider this method in the simulations of section 8.

The Wild Bootstrap is highly robust to variance misspecification since it uses within-observation sign flips of the residuals, refitting to the resulting pseudo-observations. However, this method does not extend to the GLM context since the new pseudo-observations do not generally respect the original GLM model.

In this paper we adopt the related perspective of the sign-flipping test of Hemerik et al. (2020). In this approach, rather than the residuals, the score contributions are sign-flipped. For this approach, Hemerik et al. (2020) proved a very limited robustness to variance misspecification; only for variances that are

misspecified up to a multiplicative constant. We prove that their test remains asymptotically valid under any unknown misspecification of the variance.

2.2 Standardization of the test statistics: marginal second-moment exactness

The score-flipping test (Hemerik et al., 2020) that we use as our starting point, has a tendency to anticonservativeness for small sample sizes. We prove that this anticonservative tendency is always present when that test is applied to LMs with correctly specified variances, and observe it—in simulations—in general. To address this problem we develop a novel standardization approach to the test statistics which boosts its convergence to the nominal level. Uniquely, the proposed standardization is conditional on the random sign flips. This conditional standardization ensures marginal second moment exactness, a property not shared by other resampling-based methods (Winkler et al., 2014; Frossard and Renaud, 2021).

To understand the philosophy of the new standardization approach, it is helpful to revisit the way permutation tests achieve exact control of the type I error. To get such exact control, one needs to ensure that, if the null hypothesis is true, the joint distribution of the test statistics $T(gY)$, $g \in G$, is invariant under all transformations in a group G of the data Y . That is (Hemerik and Goeman, 2018, Definition 1), under the null hypothesis,

$$(T(g'Y))_{g' \in G} \stackrel{d}{=} (T(g'gY))_{g' \in G} \quad (1)$$

for all $g \in G$. A sufficient condition for this is that $Y \stackrel{d}{=} gY$ for all $g \in G$. Once (1) is satisfied, tests can be based on randomly selected transformations G' from G .

While (1) is easy to achieve in simple experimental designs, it is difficult or impossible to achieve in more complex settings such as GLMs. Even within the LM framework, most known methods achieve only asymptotic exactness; that is, the distributions of $(T(g'Y))_{g' \in G'}$ and $(T(g'gY))_{g' \in G'}$ are not identical for all $g \in G$, but merely converge to the same (multivariate normal) distribution. As a result, type I error control is only asymptotic.

When exact equality in (1) is impossible, a fallback option is to ensure marginal second-moment null-invariance (Commenges, 2003), i.e., to let the first and marginal second moments of the distributions in (1) to be equal in finite samples. This way, asymptotic arguments are only needed for mixed and higher order moments, which tend to converge faster. We achieve marginal second-moment null-invariance by standardizing each $T(gY)$ by its standard deviation, which depends on g . To the best of our knowledge, we are the first to propose such a per-transformation standardization.

3 Set-up and assumptions

Assume that we observe n independent observations $y = \{y_1, \dots, y_n\}^T$ from the exponential dispersion family, i.e., a density of the form (Agresti, 2015)

$$f(y_i; \theta_i, \phi_i) = \exp \left\{ \frac{y_i \theta_i - b(\theta_i)}{a(\phi_i)} + c(y_i, \phi_i) \right\}, \quad (2)$$

where θ_i and ϕ_i are respectively the canonical and the dispersion parameter. We will assume throughout the paper that the model chosen fulfills the usual regularity conditions (Azzalini, 1996, Chapter 3). We derive mean and variance of the observed outcome, respectively, as

$$\mu_i = E(y_i) = b'(\theta_i); \quad \text{var}(y_i) = b''(\theta_i)a(\phi_i),$$

where primes denote derivatives. Without loss of generality, let $a(\phi_i) = \phi_i$. We assume that the mean of y_i depends on observed covariates (x_i, z_i) through the following relation, written in vector-matrix form as

$$g(\mu) = \eta = X\beta + Z\gamma$$

where $\mu = (\mu_1, \dots, \mu_n)^T$ denotes the mean vector, $g(\cdot)$ is the link function, taken to operate elementwise on a vector, (X, Z) is the full rank design matrix with $\dim(X) = n \times 1$, $\dim(Z) = n \times q$, where q does not depend on n , and (β, γ) are unknown parameters. We consider β the parameter of interest, and γ as a nuisance parameter.

Also ϕ_1, \dots, ϕ_n , are all nuisance parameters; we leave the structure of the dispersion parameters completely unspecified. Dispersion parameters can have very different values due to a variety of causes, such as the omission of some covariates (Agresti, 2015). These n separate dispersion parameters cannot be consistently estimated unless they are known to satisfy some restrictions (Neyman and Scott, 1948). Thus, in practice, the analyst will use some, possibly misspecified, putative model, which imposes such restrictions. This leads to computable estimates $\hat{\phi}_1, \dots, \hat{\phi}_n$, which will generally be inconsistent, since we cannot assume the putative model to be correct.

As an example of the set-up, consider the situation that the data are generated according to a negative binomial model with a log link function and a different and unknown dispersion parameter for every observation. The analyst could estimate these dispersion parameters under the additional constraint that $\phi_i = \phi$ for all i , or even that $\phi_i = 1$ for all i . Alternatively, the dispersion parameters could be modeled as a function of the covariates (McCullagh and Nelder, 1989, Chapter 10). Such strategies will lead to inconsistent estimates of the dispersion parameters unless the true model happens to fulfil the chosen constraints (White, 1982).

About the model and the estimation strategy we make the following assumptions.

Assumption 1. *Let (2) be the true model which generates the data. We assume that the link function $g(\cdot)$ is correctly specified.*

This assumption is crucial since it permits consistent estimation of the regression parameters β, γ , though not of ϕ_1, \dots, ϕ_n . Indeed, the incorrect estimation of the dispersion parameters ϕ_i does not affect the consistency of the estimates of the regression coefficients, as long as the mean (and hence the link function) are well-specified (Agresti, 2015). The corresponding standard errors, however, are no longer reliable.

Now, let $\ell(\beta, \gamma)$ denote the log-likelihood function, from which we can derive the score vector with elements

$$\frac{\partial \ell(\beta, \gamma)}{\partial \beta} = s_\beta = X^T D V^{-1} (y - \mu); \quad \frac{\partial \ell(\beta, \gamma)}{\partial \gamma} = s_\gamma = Z^T D V^{-1} (y - \mu),$$

where, in a compact matrix form, we have

$$D = \text{diag} \left\{ \frac{\partial \mu_i}{\partial \eta_i} \right\}; \quad V = \text{diag} \{ \text{var}(y_i) \}.$$

Taking the derivatives of the score vector we obtain the Fisher information matrix

$$J = \begin{pmatrix} J_{\beta, \beta} & J_{\beta, \gamma} \\ J_{\gamma, \beta} & J_{\gamma, \gamma} \end{pmatrix} = \begin{pmatrix} X^T W X & X^T W Z \\ Z^T W X & Z^T W Z \end{pmatrix}$$

where $W = D V^{-1} D$. Denote by d_i, v_i , and w_i the i -th diagonal elements of D, V and W , respectively. Note that all the matrices defined here and in the rest of the paper depend on n . We will suppress this dependence in the notation. The reader may assume that all quantities depend on n , except when explicitly stated otherwise.

Secondly, we assume that the estimates of ϕ_1, ϕ_2, \dots , though not consistent, will converge, as stated more formally below.

Assumption 2. For $i = 1, 2, \dots$, we have $\hat{\phi}_i - \tilde{\phi}_i \rightarrow 0$ in probability as $n \rightarrow \infty$, where $K_1 < \tilde{\phi}_i \leq K_2$, and K_1, K_2 are strictly positive constants not depending on n .

It is known (Huber, 1967) that the use of maximum likelihood estimation in a misspecified model leads under minimal conditions to a well-defined limit of the estimator. In this sense, it is possible to define a ‘‘limit’’ density, even in a misspecified model, which is intended to be the closest to the true density which generates the data in terms of Kullback-Leibler distance (White, 1982). This implies that Assumption 1 generally holds in situations where we estimate ϕ_1, \dots, ϕ_n using a restricted model.

We place a tilde symbol on all the quantities obtained when plugging the limits $\tilde{\phi}_1, \dots, \tilde{\phi}_n$ into the model. In particular, let \tilde{V}, \tilde{W} be the variance and weight matrices obtained by fixing $\phi_i = \tilde{\phi}_i$.

We further assume the following assumption related to the Fisher information given by the quantities $\tilde{\phi}_1, \dots, \tilde{\phi}_n$. This is a standard condition in regular models (Van der Vaart, 1998).

Assumption 3. Let \tilde{J} be the Fisher information matrix for $W = \tilde{W}$. The $\lim_{n \rightarrow \infty} n^{-1} \tilde{J}_{\beta, \beta}$ converges to some positive constant.

Further, the following mild condition is needed to apply the multivariate central limit theorem (Billingsley, 1986, Chapter 5) within the framework of Hemerik et al. (2020), that we will apply. It is needed to avoid pathological cases, such as vanishing or dominating observations.

Assumption 4. *Let*

$$\tilde{\nu}_{i,\beta} = \frac{(y_i - \hat{\mu}_i)x_i d_i}{\tilde{v}_i}; \quad \tilde{\nu}_{i,\gamma} = \frac{(y_i - \hat{\mu}_i)z_i d_i}{\tilde{v}_i};$$

define the element-wise score contribution and

$$\tilde{\nu}_{i,\beta}^* = \tilde{\nu}_{i,\beta} - \tilde{J}_{\beta,\gamma} \tilde{J}_{\gamma,\gamma}^{-1} \tilde{\nu}_{i,\gamma}.$$

We require, for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[(\tilde{\nu}_{i,\beta}^*)^2 \mathbf{1}_{\{|\tilde{\nu}_{i,\beta}^*|/\sqrt{n} > \epsilon\}}] \rightarrow 0$$

where $\mathbf{1}_{\{\dots\}}$ is the indicator function, and $n^{-1} \sum_{i=1}^n \text{var}(\tilde{\nu}_{i,\beta}^*)$ to converge to some positive constant.

When the model is correctly specified the last condition implies Assumption 3. In case of misspecified models we do not have this relation, since the so-called “information identity” does not hold (Azzalini, 1996).

Throughout the paper, we will make a distinction between two situations: the situation that the model is correctly specified, i.e. $\tilde{W} = W$, and the putative model is true; or the general case that $\tilde{W} \neq W$.

4 Sign-flipping effective score test

In the model described in the previous section we are interested in testing the following null hypothesis

$$H_0 : \beta = \beta_0 \mid (\gamma, \phi_1, \dots, \phi_n) \in \Gamma \times \Phi \times \dots \times \Phi \quad (3)$$

against a one or two-sided alternative, where $\Gamma \subseteq \mathbb{R}^q$ and $\Phi \subseteq (0, \infty)$. The formulation of the null hypothesis above makes explicit that only the target parameter β is fixed under H_0 , but the nuisance parameters are unconstrained.

Since ϕ_1, \dots, ϕ_n , as remarked, are difficult to estimate, we require a test that is robust to misspecification of these parameters. Hemerik et al. (2020) proposed a general semi-parametric test that has robustness to misspecification of the variance by a constant, and presented promising simulation results suggesting more general robustness properties. We will start from this test.

Hemerik et al. (2020) used the effective score for β as the test statistic. In the model (2) the effective score is

$$S = s_\beta - J_{\beta,\gamma} J_{\gamma,\gamma}^{-1} s_\gamma.$$

In the context of generalized linear models, the statistic can be written as

$$S = n^{-1/2} X^T W^{1/2} (I - H) V^{-1/2} (y - \hat{\mu})$$

where

$$H = W^{1/2} Z (Z^T W Z)^{-1} Z^T W^{1/2} \quad (4)$$

is the hat matrix, and $\hat{\mu}$ is the vector of fitted values of the model under the null hypothesis. Note that, since S is an inner product of the two n -vectors $n^{-1/2} V^{-1/2} (I - H) W^{1/2} X$ and $y - \hat{\mu}$, it can be written as a sum of n terms, which we call the effective score contributions.

To calculate the critical value, Hemerik et al. (2020) proposed to sign-flip the effective score contributions, randomly multiplying each score contribution by -1 or 1 . In matrix notation these sign flips can be represented by a random diagonal matrix \mathcal{F} of dimension n , whose non-zero elements are independent random variables that take values -1 and 1 with equal probability. Consequently, the effective sign-flip score statistic for a given flip matrix $\mathcal{F} = F$ is defined as

$$S(F) = n^{-1/2} X^T W^{1/2} (I - H) V^{-1/2} F (y - \hat{\mu}). \quad (5)$$

Note that for $\mathcal{F} = \mathbf{I}$ (the identity matrix) we recover the observed effective score.

An asymptotic α -level test is then derived as follows. First, Hemerik et al. (2020) prove (asymptotic) invariance of the first two moments of the test statistic under the action of \mathcal{F} , i.e. that $E\{S(\mathcal{F})\} - E\{S(\mathbf{I})\} = 0$ and $\text{var}\{S(\mathcal{F})\} - \text{var}\{S(\mathbf{I})\} \rightarrow 0$ as $n \rightarrow \infty$. Next, they apply the Lindberg-Feller multivariate central limit theorem to show that, for independently drawn flip matrices F_2, \dots, F_g , where g does not depend on n , the vector $S(\mathbf{I}), S(F_2), \dots, S(F_g)$ converges to a vector of independent and identically distributed random variables. Lemma 1 of Hemerik et al. (2020) then allows to obtain an asymptotic α -level test for the null hypothesis (3) against a one or two-sided alternative. Abbreviate $S_i = S(F_i)$ and denote the corresponding sorted values as $S_{(1)} \leq \dots \leq S_{(g)}$. Without loss of generality, consider testing (3) against $H_1 : \beta > \beta_0$. The test of Hemerik et al. (2020) rejects the null hypothesis if

$$S_1 > S_{(\lceil (1-\alpha)g \rceil)} \quad (6)$$

where $\lceil \cdot \rceil$ represents the ceiling function. Analogous procedures for $H_1 : \beta < \beta_0$ or $\beta \neq \beta_0$ are straightforward.

The definition of the test involves W , which is unknown. In practice we only have an estimate of W available, which converges to \tilde{W} by Assumption 1. For the theoretical results of the remainder of this paper, we will treat \tilde{W} , though not W , as known. To motivate this, we note that any error terms relating to estimation of \tilde{W} are of lower asymptotic order with respect to the parameter estimation (Barndorff-Nielsen and Cox, 1994; Pace and Salvan, 1997).

In the special case of the linear model, the effective sign-flip score test (5) simplifies to $S(F) = \sigma n^{-1/2} X^T (I - H) V^{-1/2} F (y - \hat{\mu}) = \sigma n^{-1/2} X^T (I - H) F (I - H) y$ (where (4) reduces to $H = Z (Z^T Z)^{-1} Z^T$). In that case, it

relates closely to methods that resample the residuals, summarized in the excellent reviews of Winkler et al. (2014) and Frossard and Renaud (2021). Kennedy and Cade (1996), Freedman and Lane (1983) and Dekker et al. (2007) use the same test statistic, while that of the Still-White procedure and of Draper and Stoneman (1966) is the so-called “basic score” of Hemerik et al. (2020): $\sigma n^{-1/2} X^T F(I - H)y$. Ter Braak (1992), in contrast, makes permutations of the residuals estimated under the full model. In all these methods, the null distribution is derived by permuting (or sign-flipping) the residuals and, subsequently refitting the linear model. This re-fitting approach can not be extended to GLMs, since the pseudo-responses does not retain crucial characteristics of the original response (e.g., count data may not be integer after permuting the residuals, or binomial responses are not between 0 and 1 anymore). This invalidates the required refitting step of the methods, which would also be computationally costly in GLMs. In contrast, Hemerik et al. (2020) prove that their effective score statistic (5) has asymptotically the same distribution (under the null hypothesis) both for observed and flipped test statistics, hence avoiding the refitting and permitting easy extension to the GLM context.

5 Standardized sign-flip score test

We restate the result of Hemerik et al. (2020) concerning the validity of their test below for the special case of GLMs. Our alternative proof, given in the Appendix with all other proofs, emphasizes the importance of asymptotic null-invariance.

Theorem 1. *Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumption 1-4 hold. For $n \rightarrow \infty$, the test that rejects H_0 if (6) holds is an asymptotically α level test.*

If null-invariance holds only asymptotically, the distributions of $S(\mathbf{I})$ and $S(\mathcal{F})$ are not identical, but both converge to the same distribution. It is natural to suppose that the more aspects of the distributions of $S(\mathbf{I})$ and $S(\mathcal{F})$ are equal to each other in all finite samples, the closer the resulting test is to an exact test.

The test of Hemerik et al. (2020) does not have the second-moment null-invariance property (Section 2.2): though the first moments of $S(\mathbf{I})$ and $S(\mathcal{F})$ are equal in finite samples, the variances are not. In fact, we prove in Theorem 2 below that for finite samples $S(\mathbf{I})$ always has a larger variance than $S(\mathcal{F})$, at least in the linear model with correctly specified variance. As a consequence, the test shows a tendency to anti-conservativeness in finite samples. Since the variance of the observed test statistic is always larger than its flipped counterpart, extreme values in $S(\mathbf{I})$ are more probable than in the reference distribution of $S(\mathcal{F})$. Consequently, the test tends to reject the null hypothesis too often.

Proposition 1. *Consider a normal regression model with identity link. Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumption*

1-4 hold. For finite sample size, the effective sign-flip score statistic defined as in (5) has $\text{var}\{S(\mathbf{I})\} > \text{var}\{S(\mathcal{F})\}$.

Proposition 1 is formulated for the linear model only. In the non-linear and/or variance-misspecified case $\text{var}\{S(\mathbf{I})\}$ and $\text{var}\{S(\mathcal{F})\}$ are also unequal in general. In that case the term in the asymptotic expansion that makes $\text{var}\{S(\mathbf{I})\} > \text{var}\{S(\mathcal{F})\}$ in the linear model is also present. However, in that case it is not the only asymptotic term, so it is difficult to make a fully general statement on anti-conservativeness. However, also in such models we see a tendency to small-sample anti-conservativeness in all simulations, as was also observed by Hemerik et al. (2020).

The concept of second-moment null-invariance suggests that we can improve level accuracy of the test if we can modify the flipped scores to have equal variances. The following result allows for such a procedure. It provides an expression for the variance of the flipped score, depending on the sign flip that has been applied.

Lemma 1. *The variance of the sign-flipped score, as depending on F , is*

$$\text{var}\{S(F)\} = n^{-1}X^TW^{1/2}(I - H)F(I - H)F(I - H)W^{1/2}X + o_p(1). \quad (7)$$

These variances can be estimated by plugging in $\hat{\gamma}$ and \tilde{W} . By dividing the flipped scores by their standard deviations, we obtain what we call the standardized sign-flip score statistics,

$$S^*(F) = S(F)/\text{var}\{S(F)\}^{1/2}. \quad (8)$$

We use the statistics $S^*(F)$ in the same way as the original test uses the statistics $S(F)$. The estimate of $\text{var}\{S(F)\}$ can be calculated for each F in linear time in n , as we show in Lemma 6 in the Appendix.

Analogously with the test defined in (6), consider without loss of generality testing (3) against $H_1 : \beta > \beta_0$. Take the test that rejects the null hypothesis if

$$S_1^* > S_{(\lceil(1-\alpha)g\rceil)}^*. \quad (9)$$

The following proposition establishes the null invariance and validity of the test.

Proposition 2. *Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumptions 1-4 hold. The standardized sign-flip score statistic is finite sample second-moment null-invariant. The test is asymptotically exact.*

The proposed standardization conditional on the sign flip matrix F has a similar purpose as the refitting done by methods designed for the linear model, as described in Section 2.2. It addresses both issues in that standardization than make such methods difficult to generalize to GLMs. First, there is no refitting to incongruent pseudo-observations; second, the computational effort of standardization is limited.

6 Robustness of the test

Hemerik et al. (2020) showed promising simulation results suggesting robustness against variance misspecification, but a formal proof was given only for the special case that all variances were misspecified by the same multiplicative constant. In this section we provide a formal proof of general robustness of the test to variance misspecification, both for the novel test derived in the previous section and for the original effective score test of Hemerik et al. (2020). This is the most important result of this paper.

Incorrect specification of the model variance means that

$$E \{ (y - \mu)(y - \mu)^T \} = V \neq \tilde{V},$$

where, as we recall V is the true covariance of the outcomes, and \tilde{V} is the limit of the estimates under the chosen estimation procedure. We first revisit the simple case considered by Hemerik et al. (2020), in which the variance is misspecified by a multiplicative constant. In this case, the properties of the standardized sign-flip score test are not affected by this misspecification. In particular, we have second-moment null-invariance.

Proposition 3. *Assume that Assumptions 1-4 hold. If the variances are misspecified by any finite constant $c > 0$, that is, $V = c\tilde{V}$, the standardized sign-flip score statistic is second-moment null-invariant. The test that rejects H_0 if (9) holds is asymptotically exact.*

The proposition uses the property that the test is invariant to multiplication by a constant. This result is relevant, for instance, in models with a common unknown dispersion parameter, such as normal regression model. In such cases the test can be performed, and retains second-moment null-invariance, without the need to estimate the common dispersion parameter. We can save ourselves the effort of estimating it, taking simply $\hat{\phi} = 1$. This special situation coincides with the standard parametric framework based on the quasi-likelihood approach.

The main result of this paper concerns the case of general variance misspecification. The following theorem, which is a strong improvement over the properties shown by Hemerik et al. (2020), shows that we can still get an asymptotic exact test.

Theorem 2. *Assume that the variances are misspecified, that is, $V \neq \tilde{V}$ and that Assumptions 1-4 hold. For $n \rightarrow \infty$, the standardized sign-flip score statistic is asymptotically second-moment null-invariant. The test that rejects H_0 if (9) holds is an asymptotically α level test.*

The following corollary enlarges the robustness property to the effective sign-flip score test, as an immediate consequence from the proof of the preceding theorem.

Corollary 1. *Assume that the variances are misspecified, that is, $V \neq \tilde{V}$, and that Assumptions 1-4 hold. For $n \rightarrow \infty$, the effective sign-flip score statistic is asymptotically second-moment null-invariant. The test that rejects H_0 if (6) holds is an asymptotically α level test.*

7 Multivariate test

Until now we have considered hypotheses about a single parameter $\beta \in \mathbb{R}$. We now generalize the test of the previous section to $\beta \in \mathbb{R}^d$, $d < n - q$. We consider a standard asymptotic setting where d is fixed while n increases. The null hypothesis of interest is now given by:

$$H_0 : \beta = \beta_0 \in \mathbb{R}^d \mid (\gamma, \phi_1, \dots, \phi_n) \in \Gamma \times \Phi \times \dots \times \Phi,$$

where $\Gamma \subseteq \mathbb{R}^q$ and $\Phi \subseteq (0, \infty)$, which reduces to the null hypothesis (3) if $d = 1$.

For this multivariate setting, we have to generalize the assumptions of Section 3. Assumption 2 remains unchanged while Assumptions 3 and 4 are replaced by their multivariate counterparts, respectively,

Assumption 5. *The $\lim_{n \rightarrow \infty} n^{-1} \tilde{J}_{\beta, \beta}$ converges to a positive definite matrix.*

Assumption 6. *We require, for all $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[\|\tilde{\nu}_{i, \beta}^*\|^2 \mathbf{1}_{\{\|\tilde{\nu}_{i, \beta}^*\|/\sqrt{n} > \epsilon\}}] \rightarrow \mathbf{0},$$

and $n^{-1} \sum_{i=1}^n \text{var}(\tilde{\nu}_{i, \beta}^*)$ to converge to a positive definite matrix, where $\|\cdot\|$ denotes the ℓ^2 norm and $\mathbf{0}$ is a d -dimensional zero vector.

Hemerik et al. (2020, Section 4) derived a generalization of the sign-flipped effective score test statistic, as follows. Noting that the effective score is now a d -dimensional vector $S(F) = (S^1(F), \dots, S^d(F))^T$, an asymptotically exact α -level test can be constructed by using the idea of the nonparametric combination methodology (Pesarin, 2001). The test statistic takes the form

$$T(F) = \{S(F)\}^T M \{S(F)\}$$

where M is any non-zero matrix. Usually M is chosen to be a symmetric matrix, and in general this choice influences the distribution of the power between the alternatives (see Hemerik et al. (2020) for details). A common choice for M is the inverse of an estimate of the effective Fisher information of β , if available.

As in Section 5, we can improve control of type I error by standardizing the score vector. Indeed, the same reasoning of Theorem 2 applies, showing that the sign-flip standardized score vector is finite sample second-moment null-invariant. The definition of the test is analogous by noting that now

$$\text{var}\{S(F)\} = n^{-1} X^T W^{1/2} (I - H) F (I - H) F (I - H) W^{1/2} X + o_p(1)$$

is a $d \times d$ matrix and hence the standardized score

$$S^*(F) = S(F) / \text{var}\{S(F)\}^{1/2}$$

is a d -dimensional vector. We can therefore define the test statistic as

$$T^*(F) = \{S^*(F)\}^T M \{S^*(F)\}.$$

Assume we observe values T_1^*, \dots, T_g^* , where $T_1^* = T^*(\mathbf{I})$ is the observed test statistic, while the sorted values are $T_{(1)}^* \leq \dots \leq T_{(g)}^*$. Consider the test that rejects the null hypothesis if

$$T_1^* > T_{(\lceil(1-\alpha)g\rceil)}^*. \quad (10)$$

The following Proposition states that the test is second-moment null-invariant and asymptotically exact.

Proposition 4. *Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumptions 1, 2, 5 and 6 hold. The d -dimensional standardized sign-flip score vector is finite sample second-moment null-invariant. The test that rejects H_0 if (10) holds is an asymptotically α level test.*

Finally, the following theorem shows that the robustness properties of Theorem 2 are inherited by this multivariate extension.

Theorem 3. *Assume that the variances are misspecified, that is, $V \neq \tilde{V}$ and that Assumptions 1, 2, 5, 6 hold. For $n \rightarrow \infty$, the d -dimensional standardized sign-flip score vector is asymptotically second-moment null-invariant. The test that rejects H_0 if (10) holds is an asymptotically α level test.*

8 Simulation study

We explore six different settings to compare empirically the type I error control of the usual parametric approach (by considering the Wald test), the effective and standardized flip-score tests, the Wald test based on the use of the sandwich estimator of the variance to correct for variance misspecification, the GAMLSS, modeling the variance as a function of the full model, and the HulC method, inverting the estimated confidence interval to perform hypothesis testing. A total of 5000 simulations have been carried out for each setting. The covariates have been drawn from a multivariate normal distribution, with $X \in \mathbb{R}$ and $Z \in \mathbb{R}^3$. The three nuisance covariates have correlation with the target variable equal to (0.5, 0.1, 0.1), while the true parameter is set to $\beta = 0$. The null hypothesis considered is $H_0 : \beta = 0$ against a two-sided alternative, with a nominal significance level $\alpha = 0.05$. Different sample sizes are considered (25, 50, 100, 200, 500, 1000). The sign-flip tests are performed through the R library `flipscores`. Note that for the HulC method we could not perform simulations for the smallest sample size, since after the sample splitting there were not enough observations (only four) to evaluate the model in the sub-samples.

The top two plots of Figure 1 are settings with correctly specified models, respectively a Poisson and a Logistic regression models. The middle two plots represent normal models with neglected heteroscedasticity, which depends either on a nuisance covariate or on the tested variable, i.e. respectively $\text{var}(y_i) = 4z_i^2$ and $\text{var}(y_i) = 4x_i^2$. In the bottom-left plot of Figure 1 a Poisson model was fitted when the true distribution was negative binomial with dispersion parameter

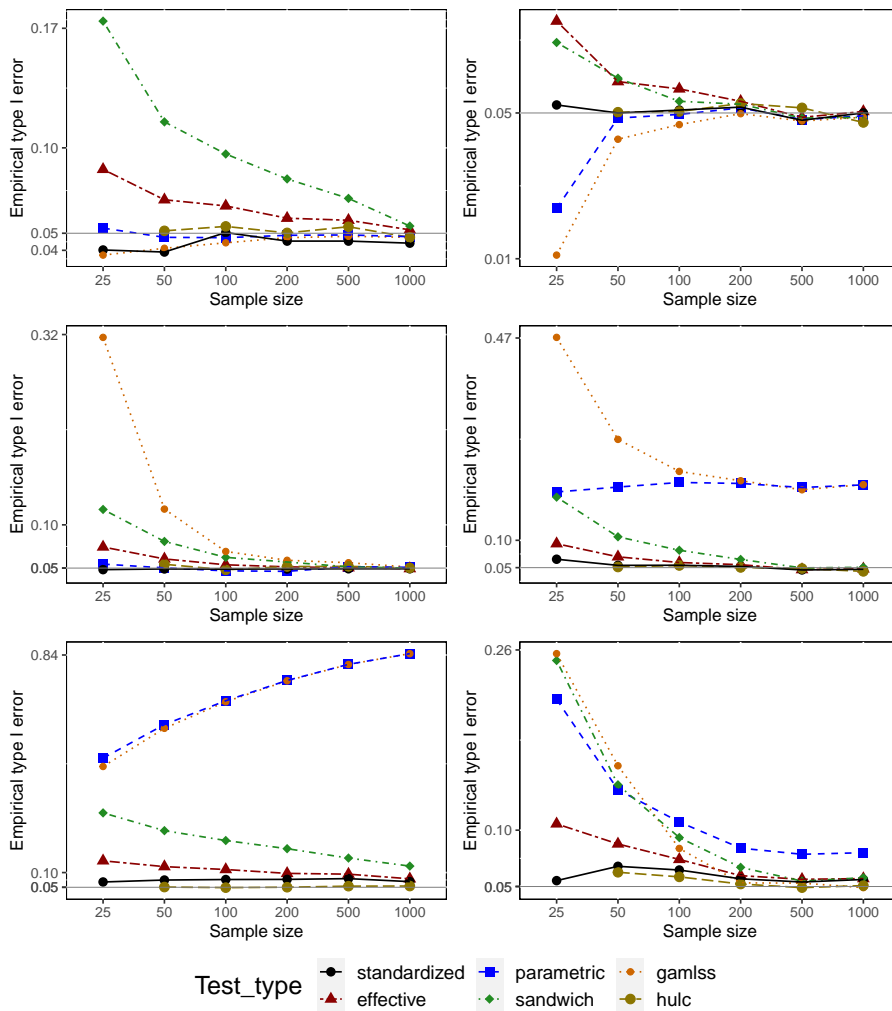


Figure 1: Type I error control comparison. **Top-left:** correct Poisson model. **Top-right:** correct Logistic model. **Middle-left:** Normal with nuisance heteroscedasticity. **Middle-right:** Normal with target heteroscedasticity. **Bottom-left:** false Poisson model. **Bottom-right:** two groups Negative-binomial.

$\phi = 1$, that is, additional heteroscedasticity relative to the Poisson model that depends on the mean. The bottom-right plot displays the results of a two-sample test fitting a negative binomial. A common dispersion parameter was assumed for the fitted negative binomial model, but the two groups are unbalanced (proportion equal to 2/3 and 1/3), and generated from two distributions with different dispersion parameter (0.4 and 1).

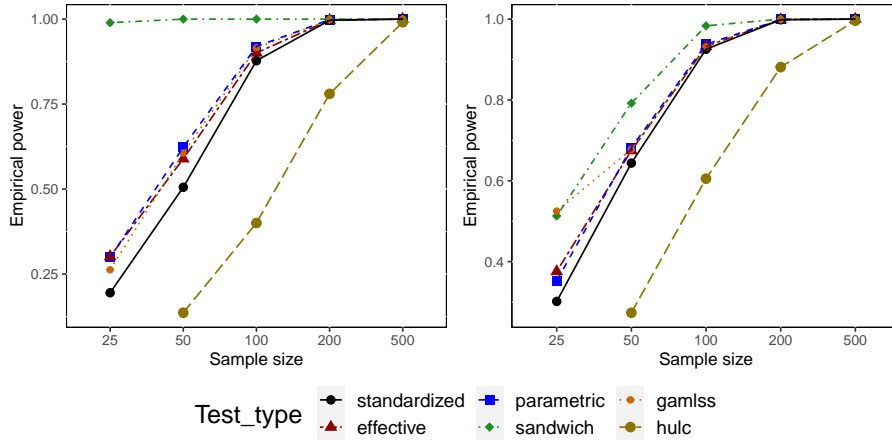


Figure 2: Power comparison. **Left:** correct Poisson model. **Right:** correct Normal model.

In the two top plots we observe that, in case of correctly specified models, the standardized test is close to the nominal level, even with $n = 25$, improving the asymptotic convergence of the effective test. The parametric test shows few rejections with small sample size for the logistic model, due to a poor approximation of the likelihood. The sandwich shows slow convergence in both cases, which is unsatisfactory since we are dealing with a well-specified model. The GAMLSS shows a low convergence, explained by the greater number of parameters involved. The HulC method is always appreciably close to the nominal level. The middle-left plot shows a similar behavior.

The last three plots show the failure of the parametric test in presence of some forms of variance misspecification, where the rejection fraction converges to a level far from the nominal for increasing sample size. The standardized test outperforms its competitors in all cases, being closest to the nominal level. In particular, the improvement over the effective test and further over the sandwich test is clear. Further, the GAMLSS surprisingly shows a failure in the middle-right plot, while in the bottom-left the reason is due to the impossibility of modeling the variance adopting a Poisson model. In the last plot it shows a slower convergence compared to the sign-flip tests. On the other side the HulC method is always close to the nominal level, with the best performance in the two bottom plots.

In the final two plots it is remarkable to see that the standardized test seems to converge to the effective test faster than going to the nominal level, leading to an initial worsening of the true level, with a recovery for larger sample size.

Finally, Figure 2 contains an evaluation of the power of the tests with two well-specified models, respectively a Poisson and a Normal model, with true parameter equal respectively to 0.3 and 1. The results for the sandwich estimator are given for completeness, although they are not comparable for small

sample size, since that method has no control of the type I error. We see that the improvement of type I error control of the standardized test with respect to the effective naturally costs some power. Analogously, we see some power loss also with respect to the parametric model and the GAMLSS as expected, but this difference is remarkably small. On the other side, the loss in power of the HulC method is remarkable, much higher with respect to the standardized test.

9 Real data analysis: RNA-Seq data

In RNA sequencing data (Love et al., 2014) a common aim is to find genes that are differentially expressed across a group of units. The usual analysis adopts a negative binomial regression model, since the observed target variables are counts, and overdispersion relative to the Poisson distribution is standard. However, the variance model generally assumes a fixed mean-variance structure with a common dispersion parameter among the groups of interest, which can be problematic, as we will show.

From the Cancer Genome Atlas (TCGA) (Tomczak et al., 2015), we have taken the TCGA-LIHC dataset of Liver Hepatocellular Carcinoma (HCC) (Erickson et al., 2016). The TCGA-LIHC consists of 20,119 genes for 344 patients with a primary tumor. We performed a very limited pre-filtering, deleting only the genes with zero total count. The target covariate is the pathological stage of the tumor. We treat it as a binary variable, splitting it between first pathological stage versus all higher stages. A total of 170 patients have a first pathological stage of the tumor, while 174 patients have higher stages. We further included two covariates in the fitted model: gender and age.

The state-of-the-art method for the analysis of these data, DESeq2 (Love et al., 2014), uses a negative binomial model with a dispersion parameter that is allowed to differ between genes, but does not depend on covariates. In particular, it does not depend on pathological stage. We tested this assumption for each gene using a GAMLSS model. The null hypothesis that dispersion did not depend on tumor stage was rejected for 5,967 out of 20,119 genes at the unadjusted 5% level. Since we would expect approximately 1,000 rejections if the DESeq2 model fits well, this gives clear indication of lack of fit of that model, at least in some of the genes. Based on the simulations in Section 8 we would, therefore, expect DESeq2 to be anti-conservative for these data.

Next, we fitted both the Poisson regression model and the negative binomial regression to each gene, and applied the Standardized sign-flip test using 5 000 flips (the default choice). We adjusted for multiple testing using the Benjamini-Hochberg correction (Benjamini and Hochberg, 1995) at $\alpha = 0.05$. We compared the results with those of DESeq2 (Love et al., 2014). This procedure estimates the nuisance parameters with an approach that shares information across genes. After obtaining a p-value for each gene, DESeq2 also applies the Benjamini-Hochberg correction at $\alpha = 0.05$ internally.

DESeq2 outputs a total of 1 450 NA-values, due to its automatic pre-filtering: 1 059 genes were filtered out due to detected outliers and 391 because of low

counts. This pre-filtering is meant to increase the power of the method, since it reduces the multiple testing burden, while removing genes for which the method has low power anyway. In contrast, the Negative Binomial regression gave 88 NA-values due to lack of convergence, while the Poisson regression gives no errors.

Table 1 gives the raw number of rejections for DESeq2 and the sign-flip tests in the second column, while the third column shows the number of rejections considering only the genes where no methods returned NA-values.

Method	No. of rejections	Filtered No. of rejections
DESeq2	3 360	3 358
Poisson sign-flip	5 109	5 032
Negative binomial sign-flip	4 833	4 765

Table 1: Number of rejections for three methods

We found that for this dataset, the standardized sign-flip test for the Poisson and Negative Binomial regression perform similarly, in fact they share the 96.1% of the conclusions. This confirms that our test is not so sensitive to (wrong) assumptions about the variance function. In contrast, DESeq2 obtains fewer rejections than both sign flip methods, despite the risk of anti-conservativeness due to the misspecification of the model.

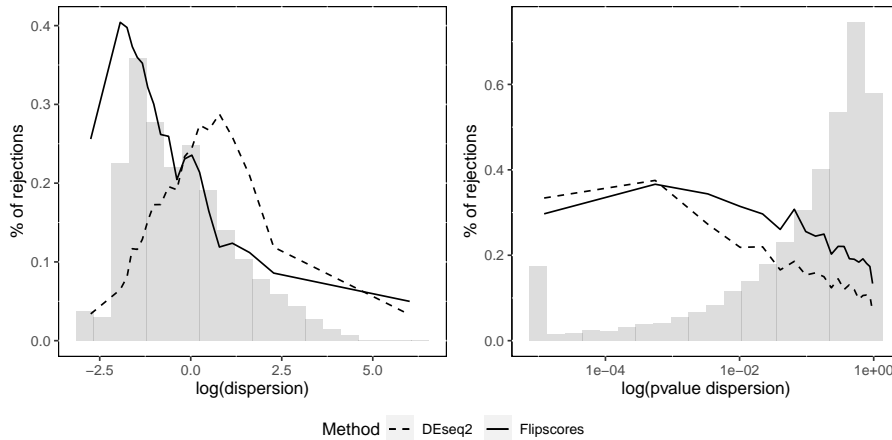


Figure 3: Percentage of rejections **Left:** compared to the dispersion parameter. **Right:** compared to the p-value of the test for equality of dispersion between groups.

Figure 3 shows the proportion of rejections as function of the maximum likelihood estimates of the dispersion parameters, on the left, and as function of the p -value of the test of equality of the dispersion parameter between the two groups. Further, the overlaid histogram represents the frequencies of the

values on the x axis. Since the two sign-flip tests have very similar results, we plot only the Negative Binomial sign-flip test against DESeq2. In the left-hand panel we observe higher power of the sign-flip test when the dispersion parameter is low, which is true for most genes. In the right-hand panel, we see a high number of rejections of DESeq2 when there is evidence for lack of equality of the dispersion parameter between groups, which can be explained from its lack of type I error control when this key assumption does not hold. When there is no such evidence, we observe higher power for the sign-flip test.

10 Discussion

Variance misspecification can dramatically reduce the quality of inference in any chosen model. We have seen that it is especially true for type I error control in hypothesis testing. In this paper we have developed a method which can be applied in the broad class of Generalized Linear Models, where it is often difficult to properly check the assumption about the variance structure. Proper variance modeling in generalized linear models is crucial in important application areas such as RNA sequencing.

We have derived a novel sign-flipping test with two important properties. First, if the model is correctly specified, the test is marginally second-moment null-invariant. As a consequence, it converges to the nominal level extremely fast, with comparable or sometimes even better control for small sample size than the parametric test. Second, if the variance model is misspecified in any arbitrary way, the new test is still asymptotically correct under minimal assumptions. Simulations show a huge improvement in small sample performance over tests based on the sandwich estimator or GAMLSS models. On the other side, the HulC method is comparable for the type I error control in misspecified problems, but only our method is applicable for the smallest sample sizes and, further, the consequences of sample splitting are visible, in the sense that we observe a relevant loss of power of HulC compared to our proposed method.

We have emphasized the value of marginal second-moment null-invariance, even though this does not imply full marginal invariance and consequently does not result in finite sample exactness. However, we note that in GLMs, unlike in the linear model, parametric tests are also unable to achieve exact control of the type I error.

References

- Alan Agresti. *Foundations of Linear and Generalized Linear Models*. Wiley, New York, 2015.
- Adelchi Azzalini. *Statistical Inference based on the likelihood*. Chapman and Hall, Boca Raton, FL, 1996.

- Ole E. Barndorff-Nielsen and David R. Cox. *Inference and Asymptotics*. Chapman and Hall, Boca Raton, FL, 1994.
- Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B*, 57(1):289–300, 1995.
- Patrick Billingsley. *Probability and Measure*. Wiley, New York, 1986.
- Dennis D. Boos. On generalized score tests. *The American Statistician*, 46(4):327–333, 1992.
- Trevor S. Breusch and Adrian R. Pagan. A simple test for heteroscedasticity and random coefficient variation. *Econometrica*, 47(5):1287–1294, 1979.
- Ching-Hui Chang and Nabendu Pal. A revisit to the Behrens–Fisher problem: Comparison of five test methods. *Communications in Statistics—Simulation and Computation*, 37(6):1064–1085, 2008.
- Daniel Commenges. Transformations which preserve exchangeability and application to permutation tests. *Nonparametric statistics*, 15(2):171–185, 2003.
- R. Dennis Cook and Sanford Weisberg. Diagnostics for heteroscedasticity in regression. *Biometrika*, 70(1):1–10, 1983.
- David R Cox. Tests of separate families of hypotheses. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability*, volume 1, page 96, 1961.
- David R Cox. Further results on tests of separate families of hypotheses. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 24(2):406–424, 1962.
- Russell Davidson and Emmanuel Flachaire. The wild bootstrap, tamed at last. *Journal of Econometrics*, 146(1):162–169, 2008.
- David Dekker, David Krackhardt, and Tom AB Snijders. Sensitivity of mrqap tests to collinearity and autocorrelation conditions. *Psychometrika*, 72:563–581, 2007.
- Norman R Draper and David M Stoneman. Testing for the inclusion of variables in linear regression by a randomisation technique. *Technometrics*, 8(4):695–699, 1966.
- Friedhelm Eicker. Reducing tcb complexity for security-sensitive applications: three case studies. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, pages 59–82, 1967.

- Bradley J. Erickson, Shanah Kirk, Y. Lee, Oliver Bathe, Melissa Kearns, C. Gerdes, Kimberly Rieger-Christ, and John Lemmerman. The cancer genome atlas liver hepatocellular carcinoma collection (tcga-lihc) (version 5) [data set]. *The Cancer Imaging Archive*, 2016.
- Ronald A. Fisher. The fiducial argument in statistical inference. *Annals of Eugenics*, 6:391–398, 1935.
- Ronald A. Fisher. The asymptotic approach to Behrens’ integral with further tables for the d test of significance. *Annals of Eugenics*, 11:141–172, 1941.
- David Freedman and David Lane. A nonstochastic interpretation of reported significance levels. *Journal of Business & Economic Statistics*, 1(4):292–298, 1983.
- David A. Freedman. On the so-called “Huber sandwich estimator” and “robust standard errors”. *The American Statistician*, 60(4):299–302, 2006.
- Jaromil Frossard and Olivier Renaud. Permutation tests for regression, ANOVA, and comparison of signals: The permuco package. *Journal of Statistical Software*, 99(15):1–32, 2021. doi: 10.18637/jss.v099.i15.
- Herbert Glejser. A new test for heteroskedasticity. *Journal of the American Statistical Association*, 64(325):316–323, 1969.
- Stephen M. Goldfeld and Richard E. Quandt. Some tests for homoscedasticity. *Journal of the American Statistical Association*, 60(310):539–547, 1965.
- Jesse Hemerik and Jelle J. Goeman. Exact testing with random permutations. *TEST*, 27:811–825, 2018.
- Jesse Hemerik, Jelle J. Goeman, and Livio Finos. Robust testing in generalized linear models by sign flipping score contributions. *Journal of The Royal Statistical Society Series B-statistical Methodology*, 82(3):841–864, 2020.
- Peter J. Huber. The behavior of maximum likelihood estimates under nonstandard conditions. *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, pages 221–233, 1967.
- Peter J. Huber and Elvezio M. Ronchetti. *Robust Statistics*. Wiley, New York, 2009.
- Carlos M. Jarque and Anil K. Bera. Efficient tests for normality, homoscedasticity and serial independence of regression residuals. *Economics Letters*, 6(3):255–259, 1980.
- Goeran Kauermann and Raymond J. Carroll. The sandwich variance estimator: efficiency properties and coverage probability of confidence intervals. *Discussion Paper 189*, 2000.

- Peter E Kennedy and Brain S Cade. Randomization tests for multiple regression. *Communications in Statistics-Simulation and Computation*, 25(4):923–936, 1996.
- Seock-Ho Kim and Allan S. Cohen. On the Behrens–Fisher problem: a review. *Journal of Educational and Behavioral Statistics*, 23(4):356–377, 1998.
- Arun Kumar Kuchibhotla, Sivaraman Balakrishnan, and Larry Wasserman. The hulc: confidence regions from convex hulls. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkad134, 2023.
- J. Scott Long and Laurie H. Ervin. Using heteroscedasticity consistent standard errors in the linear regression model. *The American Statistician*, 54(3):217–224, 2000.
- Michael I. Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology*, 15(12):1–21, 2014.
- Cora J. M. Maas and Joop J. Hox. Robustness issues in multilevel regression analysis. *Statistica Neerlandica*, 58(2):127–137, 2004.
- Jan Magnus and Heinz Neudecker. *Matrix Differential Calculus with Applications in Statistics and Econometrics*. Wiley, New York, 2019.
- Peter McCullagh and John Nelder. *Generalized Linear Models*. Chapman and Hall, Boca Raton, FL, 1989.
- Jerzy Neyman and Elizabeth L. Scott. Consistent estimates based on partially consistent observations. *Econometrica*, 16:1–32, 1948.
- Luigi Pace and Alessandra Salvan. *Principles of statistical inference: from a Neo-Fisherian perspective*, volume 4. World scientific, 1997.
- Fortunato Pesarin. *Multivariate permutation tests: with applications in biostatistics*. Wiley, Chichester, 2001.
- Robert A. Rigby and D. Mikis Stasinopoulos. Generalized additive models for location, scale and shape. *Journal of the Royal Statistical Society: Series C*, 54(3):507–554, 2005.
- Justine Rochon, Matthias Gondan, and Meinhard Kieser. To test or not to test: Preliminary assessment of normality when comparing two independent samples. *BMC Medical Research Methodology*, 12(81), 2012.
- Cajo JF Ter Braak. Permutation versus bootstrap significance tests in multiple regression and anova. In *Bootstrapping and Related Techniques: Proceedings of an International Conference, Held in Trier, FRG, June 4–8, 1990*, pages 79–85. Springer, 1992.

- Katarzyna Tomczak, Patrycja Czerwińska, and Maciej Wiznerowicz. The cancer genome atlas (tcga): an immeasurable source of knowledge. *Contemporary Oncology*, 19:A68 – A77, 2015.
- Aad W. Van der Vaart. *Asymptotic Statistics*. Cambridge University Press, Cambridge, 1998.
- Halbert White. Maximum likelihood estimation of misspecified models. *Econometrica*, 50(1):1–25, 1982.
- Anderson M. Winkler, Gerard R. Ridgway, Matthew A. Webster, Stephen M. Smith, and Thomas E. Nichols. Permutation inference for the general linear model. *Neuroimage*, 92:381–397, 2014.

A Appendix: Proofs of the Theorems

We first introduce some Lemmas needed for the proofs of the Theorems.

The projection matrix for GLMs introduced in (4) is a proper projection matrix for studentized units as shown in the following Lemma, which is mentioned for instance in (Agresti, 2015, p. 136). We give a full proof for the sake of completeness, since that is omitted in most textbooks.

Lemma 2. *The fitted and observed values in a GLM are connected through the relation*

$$V^{-1/2}(\hat{\mu} - \mu) = HV^{-1/2}(y - \mu)\{1 + o_p(1)\}$$

Proof. The first-order approximation of the score function is

$$\begin{aligned} s_{\hat{\beta}} &= s_{\beta} + \mathcal{J}_{\beta, \beta}(\hat{\beta} - \beta) + o_p(1) \\ 0 &= X^T DV^{-1}(y - \mu) - X^T WX(\hat{\beta} - \beta) + o_p(1) \\ \hat{\beta} - \beta &= (X^T WX)^{-1} X^T DV^{-1}(y - \mu) + o_p(1) \end{aligned}$$

where $(o_p(1))$ is an error term asymptotically negligible. Then, by the Delta method we have

$$\begin{aligned} \hat{\mu} &= \mu + \frac{d\mu}{d\eta} \frac{d\eta}{d\beta} (\hat{\beta} - \beta) \{1 + o_p(1)\} \\ \hat{\mu} - \mu &= DX^T (\hat{\beta} - \beta) \{1 + o_p(1)\} \\ &= DX^T (X^T WX)^{-1} X^T DV^{-1}(y - \mu) \{1 + o_p(1)\} \\ &= DW^{-1/2} HW^{1/2} D^{-1}(y - \mu) \{1 + o_p(1)\} \\ &= V^{1/2} HV^{-1/2}(y - \mu) \{1 + o_p(1)\} \\ V^{-1/2}(\hat{\mu} - \mu) &= HV^{-1/2}(y - \mu) \{1 + o_p(1)\}. \end{aligned}$$

where the asymptotically negligible error term has two sources, one related to the second-order approximation of the likelihood and the other related to possible non-linearity of the link function. \square

Lemma 3. *Let C be any n -dimensional matrix and G be a nonsingular $n \times n$ matrix, then C and $G^{-1}CG$ have the same set of eigenvalues (with the same multiplicities).*

Proof. See (Magnus and Neudecker, 2019, p. 15). \square

Lemma 4. *Let C be any n -dimensional matrix and \mathbb{F} be the set of all n -dimensional flipping matrices, i.e. the set of all possible n -dimensional diagonal matrices F with elements -1 or 1 . Then*

$$\sum_{F \in \mathbb{F}} FCF = 2^n \text{diag}(C),$$

where $\text{diag}(C)$ is a diagonal matrix with the same diagonal elements of C .

Proof. First we note that the absolute value of each element of C does not change after the multiplication FCF . The sums for the off-diagonal elements contain an equal number of terms with positive and negative sign, hence their sum over all possible flips is zero, while the sign of each diagonal element is positive for each term. By noting that the total number of flips is 2^n we have the claim. \square

The following result was adapted from Huber and Ronchetti (2009), extending their theorem (Proposition 7.1, p. 156) for linear regression models to GLMs.

Lemma 5. *Assume that the regression coefficients of a generalized linear model are consistently estimated, in the sense that for every $\varepsilon > 0$, as $n \rightarrow \infty$,*

$$\max_{1 \leq i \leq n} \text{pr}(|\hat{\mu}_i - \mu_i| > \varepsilon) \rightarrow 0.$$

Then

$$\max_{1 \leq i \leq n} h_{ii} \rightarrow 0,$$

where h_{ik} is the ik -th element of the matrix H , as defined in (4).

Proof. Using Lemma 2 we have, for each i , $\hat{\mu}_i - \mu_i =$

$$h_{ii}(y_i - \mu_i) + \sum_{k \neq i} v_i^{1/2} v_k^{-1/2} h_{ik}(y_k - \mu_k) + o_p(1). \quad (11)$$

Now we give a general probability result. Let V_1 and V_2 be two independent random variables, for any $\varepsilon > 0$ we have that

$$\begin{aligned} \text{pr}(|V_1 + V_2| \geq \varepsilon) &\geq \text{pr}(V_1 \geq \varepsilon)\text{pr}(V_2 \geq 0) + \text{pr}(V_1 \leq -\varepsilon)\text{pr}(V_2 \leq 0) \\ &\geq \min \{ \text{pr}(V_1 \geq \varepsilon), \text{pr}(V_1 \leq -\varepsilon) \}. \end{aligned}$$

Noting that $(y_i - \mu_i)$ is independent from $(y_k - \mu_k)$ for each $k \neq i$ we can apply the result to expression (11) to obtain

$$\text{pr}(|\hat{\mu}_i - \mu_i| \geq \varepsilon) \geq \min \left[\text{pr} \left\{ (y_i - \mu_i) \geq \frac{\varepsilon}{h_{ii}} \right\}, \text{pr} \left\{ (y_i - \mu_i) \leq -\frac{\varepsilon}{h_{ii}} \right\} \right] = m_i^n.$$

We have $\max_{1 \leq i \leq n} \text{pr}(|\hat{\mu}_i - \mu_i| > \varepsilon) \rightarrow 0$, so $\max_{1 \leq i \leq n} m_i^n \rightarrow 0$. Since ε was arbitrary, this implies that $\max_{1 \leq i \leq n} h_{ii} \rightarrow 0$. \square

Lemma 6. *The computational cost of the standardization constant in (8) is linear in n .*

Proof. Define $W^{1/2}Z = U\Delta L^T$, that is, the singular value decomposition of $W^{1/2}Z$, where U is a semiorthogonal $n \times q$ matrix (q equal to the rank of Z), Δ a diagonal q matrix and L a $q \times q$ orthogonal matrix. Therefore (4) can be written as

$$H = W^{1/2}Z(Z^TWZ)^{-1}Z^TW^{1/2} = U\Delta L^T L\Delta U^T = UU^T.$$

Now, let $a = (I - H)W^{1/2}X$ and $A = \text{diag}(a)$. Further, let $\mathbf{1}$ be an n -dimensional vector of ones. The denominator of the standardized test statistic becomes

$$\begin{aligned} & X^TW^{1/2}(I - H)F(I - H)F(I - H)W^{1/2}X = \\ &= a^TF(I - UU^T)Fa = \\ &= a^TFIFa - a^TFUU^T Fa \\ &= a^T a - \mathbf{1}^T AFUU^T FA\mathbf{1} \\ &= a^T a - \mathbf{1}^T FAUU^T AF\mathbf{1} \\ &= a^T a - f^T CC^T f, \end{aligned}$$

where $f = F\mathbf{1}$ and $C = AU$.

Therefore, given that $a^T a$ is a constant and the computational cost of $f^T CC^T f$ is linear with n since we can write $f^T CC^T f = \sum_{j=1}^q (\sum_{i=1}^n f_i C_{ij})^2$, the result of the lemma follows. \square

A.1 Proof of Theorem 1

Theorem. *Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumption 2-4 hold. For $n \rightarrow \infty$, the test that rejects H_0 if (6) holds is an asymptotically α level test.*

Proof. By the definition of the random sign-flipping transformations is trivial to observe that the expected value of the test statistic $S(\mathcal{F})$ is zero.

By Lemma 2 we can rewrite the effective score statistic (5) as

$$S(\mathcal{F}) = n^{-1/2}X^TW^{1/2}(I - H)\mathcal{F}(I - H)V^{-1/2}(y - \mu) + o_p(1).$$

Let

$$a = (I - H)W^{1/2}X. \tag{12}$$

The variance conditional on $\mathcal{F} = F$ is

$$\begin{aligned} \text{var}\{S(F)\} &= n^{-1}a^TF(I - H)V^{-1/2}E\{(y - \mu)(y - \mu)^T\}V^{-1/2}(I - H)Fa + o_p(1) \\ &= n^{-1}a^TF(I - H)V^{-1/2}VV^{-1/2}(I - H)Fa + o_p(1) \\ &= n^{-1}a^TF(I - H)Fa + o_p(1) \end{aligned}$$

and for $\mathcal{F} = \mathbf{I}$ we have

$$\text{var}\{S(\mathbf{I})\} = n^{-1}a^T I a + o_p(1).$$

Taking the difference

$$\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\} = n^{-1}a^T \{I - F(I - H)F\} a + o_p(1),$$

we note that the first term is a quadratic form and we look at the matrix

$$I - F(I - H)F = I - FF + FHF = FHF.$$

Since H is a projection matrix and $F^{-1} = F$, Lemma 3 implies that FHF is positive semidefinite, so that asymptotically $\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\} \geq 0$.

Define

$$h_{\text{sup}} = \sup_{1 \leq i \leq n} (h_{ii}).$$

Let us make the randomness of the flips explicit. Let \mathbb{F} denote the set of all possible flipping matrices, and note that $|\mathbb{F}| = 2^n$. By Lemma 4 we have

$$\begin{aligned} E[\text{var}\{S(\mathbf{I})\} - \text{var}\{S(\mathcal{F}) \mid F\}] &= (2^n)^{-1} \sum_{F \in \mathbb{F}} [\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\}] \\ &= n^{-1} a^T \text{diag}(H)a + o_p(1) \\ &\leq n^{-1} \|a\|^2 \cdot h_{\text{sup}} + o_p(1), \end{aligned}$$

where, by Assumption 3 and Lemma 5, the limiting behavior is

$$\lim_{n \rightarrow \infty} n^{-1} \|a\|^2 \cdot h_{\text{sup}} = 0.$$

According to the law of total variance we have

$$\text{var}\{S(\mathcal{F})\} = \text{var}[E\{S(\mathcal{F}) \mid F\}] + E[\text{var}\{S(\mathcal{F}) \mid F\}].$$

We know that $E\{S(\mathcal{F}) \mid F\}$ does not depend on \mathcal{F} , which means that $\text{var}[E\{S(\mathcal{F}) \mid F\}] = 0$, so $\text{var}\{S(\mathcal{F})\} = E[\text{var}\{S(\mathcal{F}) \mid F\}]$. It follows that, marginally over F ,

$$\lim_{n \rightarrow \infty} \text{var}\{S(\mathbf{I})\} - \text{var}\{S(\mathcal{F})\} = 0.$$

Since the random sign-flipping transformations are all independent, the corresponding test statistics are all uncorrelated, i.e., for all $1 \leq l < m \leq g$ we have $\text{cov}\{S(\mathcal{F}_l), S(\mathcal{F}_m)\} = 0$. By the previous results and Assumption 4, $(S(\mathbf{I}), \dots, S(\mathcal{F}_g))^T$ converges to a multivariate normal distribution by the multivariate Lindberg-Feller central limit theorem (Van der Vaart, 1998), with mean vector $\mathbf{0}$ and covariance matrix $s^2 \mathbf{I}$, where s^2 is the limiting variance of the flipped test statistic. Finally, we use Lemma 1 of Hemerik et al. (2020) to conclude that the test that rejects when (6) holds is an asymptotic α level test. \square

A.2 Proof of Proposition 1

Proposition. *Consider a normal regression model with identity link. Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumption 2-4 hold. For finite sample size, the effective score statistic defined as in (5) has $\text{var}\{S(\mathbf{I})\} > \text{var}\{S(\mathcal{F})\}$.*

Proof. By Lemma 2 we can rewrite the effective score statistic (5) as

$$S(\mathcal{F}) = n^{-1/2} X^T W^{1/2} (I - H) \mathcal{F} (I - H) V^{-1/2} (y - \mu) + o_p(1).$$

The variance conditional on $\mathcal{F} = F$ is

$$\begin{aligned} \text{var}\{S(F)\} &= n^{-1} a^T F (I - H) V^{-1/2} E\{(y - \mu)(y - \mu)^T\} V^{-1/2} (I - H) F a \\ &= n^{-1} a^T F (I - H) V^{-1/2} V V^{-1/2} (I - H) F a \\ &= n^{-1} a^T F (I - H) F a \end{aligned}$$

and for $\mathcal{F} = \mathbf{I}$ we have

$$\text{var}\{S(\mathbf{I})\} = n^{-1} a^T I a.$$

Taking the difference

$$\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\} = n^{-1} a^T \{I - F(I - H)F\} a,$$

we note that it is a quadratic form and we consider the matrix

$$I - F(I - H)F = I - FF + FHF = FHF.$$

Since H is a projection matrix and $F^{-1} = F$, Lemma 3 implies that FHF is positive semidefinite, and therefore $\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\} \geq 0$.

We then prove the strict inequality for at least one flipping matrix F . Taking any model with an intercept, by construction

$$h_{\inf} = \inf_{1 \leq i \leq n} (h_{ii}) > 0.$$

Note that $|\mathbb{F}| = 2^n$. Since $\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\} \geq 0$ as proven above, it suffices to show that

$$\sum_{F \in \mathbb{F}} n^{-1} a^T F H F a > 0,$$

which means that we have a strictly inequality for some F . Using Lemma 4 we have

$$\sum_{F \in \mathbb{F}} F H F = 2^n \text{diag}(H).$$

Therefore

$$\begin{aligned} \sum_{F \in \mathbb{F}} n^{-1} a^T F H F a &= 2^n n^{-1} a^T \text{diag}(H) a \\ &\geq 2^n n^{-1} \|a\|^2 \cdot h_{\inf} > 0. \end{aligned}$$

□

A.3 Proof of Lemma 1

Lemma. *The variance of the sign-flipped score, as depending on F , is*

$$\text{var}\{S(F)\} = n^{-1}X^TW^{1/2}(I-H)F(I-H)F(I-H)W^{1/2}X + o_p(1).$$

Proof. Let

$$a = (I-H)W^{1/2}X.$$

The variance for a given sign-flip matrix F is

$$\begin{aligned} \text{var}\{S(F)\} &= n^{-1}a^TF(I-H)V^{-1/2}E\{(y-\mu)(y-\mu)^T\}V^{-1/2}(I-H)Fa + o_p(1) \\ &= n^{-1}a^TF(I-H)V^{-1/2}VV^{-1/2}(I-H)Fa + o_p(1) \\ &= n^{-1}a^TF(I-H)Fa + o_p(1) \end{aligned}$$

□

A.4 Proof of Proposition 2

Proposition (copy from main text). /

Proof. We observe that the expected value of the new statistic (8) is left unchanged, while the standardization makes the variance of each flipped new statistic equal to 1. The same argument given in the last part of the proof of Theorem 1 applies to deduce the asymptotic exactness of the test.

In the special case of the normal model with the identity link, we have that $\text{var}(S(F))$ does not depend on any unknown nuisance parameters. Moreover, in this model $S(\mathbf{I})$ and $S(\mathcal{F})$ are normally distributed in finite samples. Further, the w test statistics are all uncorrelated due to the random sign-flipping of the underlying summands. Hence the standardized sign-flip score statistic is second-moment exact. □

A.5 Proof of Proposition 3

Proposition (copy from main text). /

Proof. We observe that

$$V\tilde{V}^{-1} = \phi\tilde{\phi}^{-1}I = cI$$

and therefore

$$\text{var}\{S^*(F)\} = c \quad \forall F \in \mathbb{F},$$

as was to be shown. □

A.6 Proof of Theorem 2

Theorem. *Assume that the variances are misspecified, that is, $V \neq \tilde{V}$ and that Assumptions 2-4 hold. For $n \rightarrow \infty$, the standardized sign-flip score statistic is asymptotically second-moment null-invariant. The test that rejects H_0 if (9) holds is an asymptotically α level test.*

Proof. It is trivial to see that the expected value of the test statistic is not affected by this misspecification.

Let

$$B = V\tilde{V}^{-1},$$

note that it is a diagonal matrix and all its elements are finite and greater than zero by Assumption 2. We compute again the variance of (5), but now we consider the misspecification.

Let \tilde{H} and \tilde{a} be the quantities defined in (4) and (12) for $W = \tilde{W}$. The variance can be written as

$$\begin{aligned} \text{var}\{S(F)\} &= n^{-1}\tilde{a}^T F(I - \tilde{H})\tilde{V}^{-1/2} E\{(y - \mu)(y - \mu)^T\} \tilde{V}^{-1/2}(I - \tilde{H})F\tilde{a} + o_p(1) \\ &= n^{-1}\tilde{a}^T F(I - \tilde{H})B(I - \tilde{H})F\tilde{a} + o_p(1). \end{aligned}$$

Take the difference

$$\begin{aligned} \text{var}(S(\mathbf{I})) - \text{var}(S(F)) &= n^{-1}\tilde{a}^T \left\{ B - F(I - \tilde{H})B(I - \tilde{H})F \right\} \tilde{a} + o_p(1) \\ &= n^{-1}\tilde{a}^T \left[F\{\tilde{H}B + (I - \tilde{H})B\tilde{H}\}F \right] \tilde{a} + o_p(1). \end{aligned}$$

We notice that by Lemma 4

$$\sum_{F \in \mathcal{F}} F(\tilde{H}B + B\tilde{H} - \tilde{H}B\tilde{H})F = 2^n \text{diag}(\tilde{H}B + B\tilde{H} - \tilde{H}B\tilde{H})$$

where the i -th element of that diagonal matrix is

$$2\tilde{h}_{ii}b_i - \sum_{k=1}^n \tilde{h}_{ik}^2 b_k.$$

Then we have

$$\begin{aligned} E[\text{var}\{S(\mathbf{I})\} - \text{var}\{S(\mathcal{F}) \mid F\}] &= (2^n)^{-1} \sum_{F \in \mathbb{F}} [\text{var}\{S(\mathbf{I})\} - \text{var}\{S(F)\}] \\ &= n^{-1}\tilde{a}^T \text{diag}(\tilde{H}B + B\tilde{H} - \tilde{H}B\tilde{H})\tilde{a} + o_p(1) \\ &\leq n^{-1}\tilde{a}^T \text{diag}(\tilde{H}B + B\tilde{H})\tilde{a} + o_p(1). \end{aligned}$$

By Assumption 2, note that there exists two finite positive constants c_1, c_2 such that for each i -th element

$$2\tilde{h}_{ii}b_i \leq c_1 b_{sup} \tilde{h}_{ii} \leq c_2 \sup_{1 \leq i \leq n} \tilde{h}_{ii} = c_2 \tilde{h}_{sup}$$

where

$$\tilde{h}_{sup} = \sup_{1 \leq i \leq n} (\tilde{h}_{ii}) \quad b_{sup} = \sup_{1 \leq i \leq n} (b_i).$$

Therefore we can derive the upper bound

$$E[\text{var}\{S(\mathbf{I})\} - \text{var}\{S(\mathcal{F}) \mid F\}] \leq n^{-1} \|\tilde{a}\|^2 c_2 \cdot \tilde{h}_{sup}.$$

Using Assumption 3 and Lemma 5, the limiting behavior is

$$\lim_{n \rightarrow \infty} n^{-1} \|\tilde{a}\|^2 c_2 \cdot \tilde{h}_{\text{sup}} = 0.$$

Meanwhile, for two positive constants c_3, c_4 we have for each i -th element

$$-\sum_{k=1}^n \tilde{h}_{ik}^2 b_k \geq -\sum_{k=1}^n \tilde{h}_{ik}^2 c_3 b_{\text{sup}} = -c_3 b_{\text{sup}} \tilde{h}_{ii} \geq -c_4 \tilde{h}_{\text{sup}}.$$

Then using again Lemma 4 we can derive the lower bound

$$\begin{aligned} E[\text{var}\{S(\mathbf{I})\} - \text{var}\{S(\mathcal{F}) \mid F\}] &\geq n^{-1} \tilde{a}^T \text{diag}(-\tilde{H}B\tilde{H})\tilde{a} + o_p(1) \\ &\geq -n^{-1} \|\tilde{a}\|^2 c_4 \cdot \tilde{h}_{\text{sup}} + o_p(1) \end{aligned}$$

where, using Assumption 3 and Lemma 5, the limiting behavior is

$$\lim_{n \rightarrow \infty} -n^{-1} \|\tilde{a}\|^2 c_4 \cdot \tilde{h}_{\text{sup}} = 0.$$

According to the law of total variance we have

$$\text{var}\{S(\mathcal{F})\} = \text{var}[E\{S(\mathcal{F}) \mid F\}] + E[\text{var}\{S(\mathcal{F}) \mid F\}].$$

We know that $E\{S(\mathcal{F}) \mid F\}$ does not depend on \mathcal{F} , which means that $\text{var}[E\{S(\mathcal{F}) \mid F\}] = 0$, so $\text{var}\{S(\mathcal{F})\} = E[\text{var}\{S(\mathcal{F}) \mid F\}]$. It follows that, marginally over F ,

$$\lim_{n \rightarrow \infty} \text{var}\{S(I)\} - \text{var}\{S(\mathcal{F})\} = 0.$$

The same argument given in the last part of the proof of theorem 1 applies to deduce the asymptotic exactness of the test considered. \square

A.7 Proof of Proposition 4

Proposition. *Assume that the variances are correctly specified, that is, $\tilde{V} = V$, and that Assumptions 2, 5 and 6 hold. The d -dimensional standardized sign-flip score vector is finite sample second-moment null-invariant. The test that rejects H_0 if (10) holds is an asymptotically α level test.*

Proof. The proof is analogous to Theorem 3 of Hemerik et al. (2020). \square

A.8 Proof of Theorem 3

Theorem. *Assume that the variances are misspecified, that is, $V \neq \tilde{V}$ and that Assumptions 2, 5, 6 hold. For $n \rightarrow \infty$, the d -dimensional standardized sign-flip score vector is asymptotically second-moment null-invariant. The test that rejects H_0 if (10) holds is an asymptotically α level test.*

Proof. It is trivial to observe that the expected value is a d -dimensional zero vector. We focus on the variance, which now is a $d \times d$ matrix. We prove the robustness elementwise.

The proof for the diagonal elements follows from Theorem 2. Then, we focus on one covariance term of the covariance matrix of the flipped effective score vector. Let $X = (X_1, \dots, X_d)^T$ and

$$\begin{aligned}\tilde{a}_1 &= (I - \tilde{H})\tilde{W}^{1/2}X_1 \\ \tilde{a}_2 &= (I - \tilde{H})\tilde{W}^{1/2}X_2.\end{aligned}$$

When the variances are misspecified

$$\begin{aligned}\text{cov}\{S^1(F), S^2(F)\} &= \\ &= n^{-1}\tilde{a}_1^T F(I - \tilde{H})\tilde{V}^{-1/2}E\{(y - \mu)(y - \mu)^T\}\tilde{V}^{-1/2}(I - \tilde{H})F\tilde{a}_2 + o_p(1) \\ &= n^{-1}\tilde{a}_1^T F(I - \tilde{H})B(I - \tilde{H})F\tilde{a}_2 + o_p(1).\end{aligned}$$

Take the difference

$$\begin{aligned}\text{cov}\{S^1(I), S^2(I)\} - \text{cov}\{S^1(F), S^2(F)\} &= \\ &= n^{-1}\tilde{a}_1^T \left\{ B - F(I - \tilde{H})B(I - \tilde{H})F \right\} \tilde{a}_2 + o_p(1) \\ &= n^{-1}\tilde{a}_1^T \left\{ F(\tilde{H}B + B\tilde{H} - \tilde{H}B\tilde{H})F \right\} \tilde{a}_2 + o_p(1).\end{aligned}$$

From this point the asymptotic second-moment null-invariance can be derived directly from the proof of the Theorem 2, by applying the same reasoning to each covariance term, replacing Assumptions 3-4 with Assumptions 5-6. The asymptotic exactness of the test follows from the proof of Proposition 4. \square