



JENA ECONOMIC RESEARCH PAPERS



2007 – 032

Small Sample Properties of the Wilcoxon Signed Rank Test with Discontinuous and Dependent Observations

by

**Nadine Chlaß
Jens J. Krüger**

www.jenecon.de

ISSN 1864-7057

The JENA ECONOMIC RESEARCH PAPERS is a joint publication of the Friedrich-Schiller-University and the Max Planck Institute of Economics, Jena, Germany. For editorial correspondence please contact m.pasche@wiwi.uni-jena.de.

Impressum:

Friedrich-Schiller-University Jena
Carl-Zeiß-Str. 3
D-07743 Jena
www.uni-jena.de

Max-Planck-Institute of Economics
Kahlaische Str. 10
D-07745 Jena
www.econ.mpg.de

© by the author.

Small Sample Properties of the Wilcoxon Signed Rank Test with Discontinuous and Dependent Observations

by

Nadine Chlaß

Max Planck Institute of Economics

Strategic Interaction Group, Kahlaische Strasse 10, D-07745 Jena, Germany,
Tel.: +49 3641 686 624, Fax: +49 3641 686 623, E-Mail: chlass@econ.mpg.de

and

Jens J. Krüger

Friedrich-Schiller-University Jena

Department of Economics, Carl-Zeiss-Strasse 3, D-07743 Jena, Germany,
Tel.: +49 3641 943 200, Fax: +49 3641 943 202, E-Mail: jens.krueger@wiwi.uni-jena.de

July 2007

Abstract

This Monte-Carlo study investigates sensitivity of the Wilcoxon signed rank test to certain assumption violations in small samples. Emphasis is put on within-sample-dependence, between-sample dependence, and the presence of ties. Our results show that both assumption violations induce severe size distortions and entail power losses. Surprisingly, these consequences do vary substantially with other properties the data may display. Results provided are particularly relevant for experimental settings where ties and within-sample dependence are frequently observed.

JEL classification: C12, C14, C15

Keywords: Wilcoxon signed rank test, ties, dependent observations, size and power

1 Introduction

The Wilcoxon rank tests constitute widely used nonparametric tests of sample dissimilarities based on ranked differences (Feltovic 2003, p. 274). Their virtues comprise a small number of assumptions (for an overview see Hollander/Wolfe 1999, p. 46) as well as the comparability of subjects at two points in time, thus allowing for the comparison of related samples. Wilcoxon rank tests rely on the continuity of investigated variables and therefore absence of ties, and second, the independence of observations within one sample. In what is to follow we investigate the impact of exclusive as well as simultaneous violations of these assumptions on both size and power properties for the Wilcoxon signed rank test.

The Wilcoxon signed rank test is routinely applied in experimental economics where both assumptions are frequently violated. Within experimental settings, there are several threats to continuity of response variables. First, the definition range of responses is not only finite, but frequently quite narrowly restricted, for example within $[0,10]$ (Chlaß et al. 2006) or even within $[0,1]$ (Dittrich et al. 2005). Second, actually observable responses are often furthermore reduced in variety by, for example, the response being restricted in decimals, different subjects using similar decision heuristics, or different subjects simultaneously playing the dominant strategy in each of the two points (situations) in time. Equal differences (ties) may thus frequently result and require consideration. Breaking of ties, e.g. via mid ranks (Hollander/Wolfe 1999, pp. 109), solves the problem only partially, since the question how differences could have been ranked, if more continuous responses had been possible or if subjects had thought to answer in a more continuous way cannot be answered. A more reliable approach therefore seems to consider whether the frequency of ties severely questions the validity of test results and then rather to opt for a test on discrete data (e.g. a χ^2 test), putting up with some loss of information. This, however, requires more detailed knowledge on the impact of ties - knowledge we wish to provide.

Dependence of observations within experimental settings does not only originate from repeatedly measuring the same individual's responses, but also from interaction of participants or latent variables in general. Thus, data possibly remain dependent even on an aggregated level, often rendering only the entire experimental session an independent observation. Only controlling for no interaction at all between several subgroups of subjects provides a larger number of independent observations. Within those subgroups, however, the problem of aggregation remains. Therefore, a considerable amount of data and information is

lost, decreasing overall test power and interpolating temporal trends whose investigation might be of interest as well. There are, of course, empirical methods, e.g. random or mixed effect models which account for dependence within samples and provide efficient and unbiased estimates for such cases. However, these estimates do not allow for any causal statistical interpretation in for which one would need information on the true conditional (in)dependence relations. Obtaining causal information, however, is the very aim of controlled experiments. As a consequence, dependence of observations might provide a valuable indicator of failures within an experimental design and thus serves as a control mechanism that should be retained rather than being accounted for statistically. Hence, it is, from a specification point of view, valuable to assess the impact of dependence on an estimate rather than to provide for it.

Our study investigates the effects of ties and dependent observations on size and power properties of the Wilcoxon signed rank test, determining actual p-values under each of the above-mentioned assumption violations. It aims at providing decision support as to what degree of violations the test remains a reliable tool and when its outcomes cannot be trusted. We proceed by reviewing hitherto obtained results regarding the Wilcoxon test in section 2, detail our simulation design in section 3, discuss our findings in sections 4 and 5 and conclude in section 6.

2 Literature Review

Let us first give an impression of hitherto conducted studies and the insights obtained so far regarding the Wilcoxon tests in general. Starting point of this literature is the study of Hodges and Lehmann (1956) who provide an early comparison of the Wilcoxon rank sum test and some other nonparametric alternatives to the t-test, analyzing asymptotic efficiency. They establish a lower threshold for relative Pitman efficiency of 0.864 in comparing Wilcoxon and t-test. Thus, in testing against shifts, the efficiency loss of the former as compared to the latter is bounded, while efficiency gains on the contrary may be infinite. However, the same comparison shows other rank tests, e.g. by Fisher and Yates (1948) and van der Waerden (1953), to reach relative asymptotic power of unity.

Turning to Monte-Carlo studies, Tanizaki (1997) assesses and compares power properties of the Wilcoxon rank sum and other rank-based tests to the t-test under various distributional assumptions. These comprise normal, Cauchy, logistic, chi-square and uniform distributions.

The results indicate that in small samples, the Wilcoxon test has highest overall power, performs better than other rank-based tests and clearly outperforms the t-test.

Zimmerman (1998) provides a Monte Carlo comparison of the Wilcoxon rank sum test with a modification of the t-test featuring various violations of normality and different standard deviations of the samples. Though for equal standard deviations the Wilcoxon test maintains its exact size as opposed to the t-test, neither test does so for unequal dispersions. In case of heavy-tailed distributions, different standard deviations and different sample sizes the Wilcoxon test may be more severely distorted than the t-test. The study does not investigate power properties.

Zimmerman (2000) provides a detailed Monte-Carlo investigation confined to the size properties of the Wilcoxon rank sum test. Given normally distributed samples and unequal dispersions the test systematically overrejects the null hypothesis if dispersions are unequal while the t-test appears to overreject less pronouncedly. These differences are due to the fact that the Wilcoxon test tests the equality of entire distributions instead of only considering their location parameters.

Freidlin and Gastwirth (2000) compare power properties of the Wilcoxon rank sum test to a variety of competing nonparametric tests for different underlying distributions and situations with equal or unequal sample sizes. They find the Wilcoxon rank sum test to have optimal power properties given logistic distributions. Comparative performance is rather discouraging under fat-tailed slash, Cauchy and $t(2)$ distributions, while the test still demonstrates sound power for normal and double exponential distributions.

Finally, Feltovich (2003) in an extensive study assesses both size and power of the Wilcoxon rank sum test for a variety of situations and additionally provides a comparison to the robust-rank-order test of Fligner and Policello (1981). The situations investigated cover different sample sizes which may be equal or unequal, differing standard deviations and asymmetry of the underlying distributions. Furthermore compared are exact statistics and their normal approximations. A detailed assessment of size properties finds the robust-rank-order test to react more sensitively than the Wilcoxon rank sum test. For asymmetrically distributed data, both tests are found to perform poorly and tend to overreject, with an outperforming Wilcoxon rank sum test for unequal sample sizes. Power analysis is limited to the case with different central tendencies and equal dispersions of samples, no substantial differences being observed between the two tests.

The findings reviewed here find the Wilcoxon test a competitive testing method in many situations. The assumption violations investigated are limited to different dispersions and distributional asymmetry. Two further violations of assumptions required by the test having not yet received any attention remain tied observations and within-sample dependence. While for ties, remedying techniques have been proposed, such as the algorithm by Streitberg and Röhmer (1986, 1987), we show that their consequences persist nevertheless. Within-sample dependence on the other hand, has not yet received any attention at all. Considering the fact that especially in the context of experimental data, both phenomena often occur jointly, we also investigate simultaneous assumption violations.

While previous studies have focused on the Wilcoxon (rank sum) test for unrelated samples, we will concentrate on the procedure for related samples and additionally inquire to what extent the degree of sample relatedness matters.

3 Simulation Design

Our baseline scenario starts by simulating a series $\{v_t\}$ of identically and independently distributed random draws from a normal distribution with mean zero and standard deviation σ_v . The sample size is chosen by taking into account the fraction of ties $\lambda \in [0,1)$ that will be induced later on, so that a total of λn random draws of v are simulated to achieve a sample size of n . Note that λ and n are appropriately chosen so that λn is always integer. From $\{v_t\}$ a series $\{x_t\}$ is constructed via a simple stationary first-order autoregressive process AR (1) with autocorrelation parameter $\rho_x \in [0,1)$ where one period influences only a limited number of future periods. This amounts to simulating $x_t = \rho_x x_{t-1} + v_t$, allowing for a certain burn-in phase. Since we actually want to control the variance of the process $\{x_t\}$, we control for σ_x throughout the simulations and compute σ_v via $\sigma_v = \sigma_x \sqrt{1 - \rho_x^2}$.

In a subsequent step a series $\{y_t\}$ is constructed from $\{x_t\}$ using a simple linear relationship $y_t = \delta + \beta x_t + u_t$ where $\{u_t\}$ represents another series of identically and independently distributed draws from a normal distribution with mean zero and standard deviation σ_u . Within these simulations we control for the standard deviation of y , σ_y , (in addition to the standard deviation of x , σ_x) and the correlation coefficient between x and y , ρ_{xy} . This is

implemented by computing β by $\beta = \rho_{xy} \cdot \sigma_x / \sigma_y$ and σ_u by $\sigma_u = \sqrt{\sigma_y^2 - \beta^2 \sigma_x^2}$. Thus, we are able to allow for different variances of the series $\{x_t\}$ and $\{y_t\}$ as well as for different degrees of dependence between them. Since $\{y_t\}$ originates from $\{x_t\}$, the properties of the latter carry over to the former so that $\{y_t\}$ equally exhibits stationary properties with a first-order autocorrelation coefficient $\rho_y = \rho_x \rho_{xy}^2$.

The mean difference between $\{x_t\}$ and $\{y_t\}$ is controlled by choosing δ such that the value of Hedges g represents no, small, medium or large differences. Hedges g (Hedges, Olkin 1985) is a measure similar to the more widely known Cohen's d (see Cohen 1988, p.23). Taking sample size into account, it allows for a better judgement of small, medium and large effects within comparisons of test power. In our case with equal sample sizes this just amounts to making an adjustment to the overall effect size based on the sample size. Specifically, we control δ by fixing g and computing $\delta = g \cdot \sqrt{(\sigma_x^2 + \sigma_y^2)/2} / (1 - 3/(8n - 9))$, reflecting equal sample sizes for x and y .

Ties are induced in two different forms. The first variant randomly draws $(1 - \lambda)n$ values from the λn simulated values of x and y with replacement and adds them to the sample, resulting in a total sample size of n . The second variant draws a single value out of the λn simulated values of x and y and adds this value $(1 - \lambda)n$ times to the sample, resulting again in a total sample size of n . These variants are called the “many-small-ties variant” and the “one-big-tie variant”, respectively. Of course, this way of inducing ties results in a total sample of values displaying different time series properties from those originally specified by σ_x , σ_y , ρ_x and ρ_{xy} . Thus, whenever ties are present, all subsequent interpretations of our results e.g. with respect to different degrees of dependence are only suggestive and should be viewed as qualitative statements regarding relatively higher or lower degrees of dependence.

Parameters are specified as follows. We consider four different sample sizes $n \in \{10, 20, 30, 40\}$ and four different fractions of ties $\lambda \in \{0, 0.2, 0.5, 0.8\}$. The standard deviation of x is fixed at unity, $\sigma_x = 1$, and that of y takes values $\sigma_y \in \{1, 2, 3, 4\}$. For the autocorrelation of x we investigate four different values $\rho_x \in \{0, 0.2, 0.5, 0.8\}$ as we do for the correlation between x and y , $\rho_{xy} \in \{0, 0.2, 0.5, 0.8\}$. Finally, Hedges g is specified in relation

to Cohen's d as $g \in \{0, 0.2, 0.5, 0.8\}$, representing no, small, medium and large effects as suggested by Cohen (1988 p.25).

We simulate 10000 replications for each of the 4096 possible combinations of parameter values and afterwards apply the Wilcoxon signed rank test to each of the simulated samples of x and y . A two-sided test is performed by comparing the test statistic to the respective critical values for a nominal significance level of $\alpha = 0.05$ which will be used throughout this paper. Finally recorded is the fraction of rejections out of the 10000 replications for each parameter combination.

Departing from our baseline scenario introduced above, we specifically modify the underlying distribution of samples $\{v_t\}$ in order to investigate how fat tails or asymmetry of the distribution affect the test properties under assumption violations. As representative examples for distributions with large tails and asymmetry t and gamma distributions were chosen. The first case corresponds to a series being drawn from a t distribution with three degrees of freedom. A subsequent division by its standard deviation $\sqrt{3}$ and scaling by σ_v^2 , yields an error process $\{v_t\}$ with mean zero and variance σ_v^2 . Regarding parameters for the gamma distribution controlling shape a and scale b we first require the variance being fixed at unity, $\sigma_v^2 = ab^2 = 1$ and second, prespecify skewness via $\gamma = 2/\sqrt{a}$. This results in $a = (2/\gamma)^2$ and $b = \sqrt{1/a}$. A subsequent subtraction of the mean ensures a zero-mean error process $\{v_t\}$. To shed light on the impact of skewness, parameter a is set to take values $a \in \{1, 2, 3, 4\}$.

Furthermore, we investigate the impact of specific kinds of dependence running simulations with a first order moving-average, or MA(1), process $x_t = v_t + \theta v_{t-1}$ where each period influences infinitely many future periods. This specification relies on inverting the first-order autocorrelation coefficient $\rho_x = \theta/(1+\theta^2)$ for computing θ and then using this value to compute the variance of $\{v_t\}$ by $\sigma_v^2 = \sigma_x^2/(1+\theta^2)$.

Finally, all constellations mentioned above are investigated in their impact on size and power properties for the exact version of the Wilcoxon signed rank test which, in the presence of ties, is implemented using the algorithm proposed in (Streitberg/Röhmer 1986, 1987). Our questions are first, whether ties still show a residual impact and second, whether the impact of

within-sample dependence differs anyhow compared to the standard test. To summarize this chapter, figure 1 reviews all treatment combinations investigated.

distribution	normal	test type	standard	significance level	0.01	dependence	AR	ties	many small ties		
	t(3)										one large random tie
	gamma		exact		0.05		MA		one large median tie		

Figure 1: Treatment combinations investigated.

4 Simulation Results

Table 1 presents rejection frequencies for a *first scenario* constituting our baseline case. It introduces a normal distribution, an intermediary degree of dependence $\rho_{xy} = 0.5$, equal variances $\sigma_x = \sigma_y = 1$, AR(1) errors and many small ties, drawn with replacement from the simulated values. With the samples showing no difference and in absence of any assumption violation, the test maintains its *size* level quite exactly. *Power* increases with both accruing deviation from the null hypothesis g and sample size and becomes quite considerable for larger sample sizes. However, power persists on a sometimes considerably low level for both small and medium size effects, first rising above what is recommended in (Cohen 1988, p. 53) for $\{g = 0.5, n = 40\}$.

As a first assumption violation, let us now investigate the impact of discontinuity in the investigated variables entailing tied observations on both size and power properties. The first four rows with $\rho_x = 0$ and $\lambda \in \{0, 0.2, 0.5, 0.8\}$ indicate that, when within-sample-dependence is absent but *ties*¹ are present to various degrees, the test rejects a true null hypothesis more often than it should, the problem persisting with increasing sample size. For $\lambda = 0.2$, that is,

¹ The effects of both ties and within-sample dependence on *size* as mentioned here are not immediately obtained by visual inspection of the tables provided. To give an example, the tie effect on size $\Delta_t^S|_{\{\lambda=\bar{\lambda}, n=\bar{n}\}}$ is calculated

by $\Delta_t^S|_{\{\lambda=\bar{\lambda}, n=\bar{n}\}} = (a - (a - b)) - (c - (c - d))$ with $a := f_r|_{\{g=\bar{g}, n=\bar{n}, \lambda=\bar{\lambda}, \rho_x=0\}}$, $b := f_r|_{\{g=0, n=\bar{n}, \lambda=\bar{\lambda}, \rho_x=0\}}$, $c := f_r|_{\{g=\bar{g}, n=\bar{n}, \lambda=0, \rho_x=0\}}$, $d := f_r|_{\{g=0, n=\bar{n}, \lambda=0, \rho_x=0\}}$ where f_r denotes the respective rejection frequency displayed in the corresponding table. The impact of ties on size- adjusted *power* is calculated the following way: $\Delta_t^P|_{\{\lambda=\bar{\lambda}, n=\bar{n}\}} = f_r|_{\{\lambda=\bar{\lambda}, n=\bar{n}, \rho_x=0\}} - f_r|_{\{\lambda=0, n=\bar{n}, \rho_x=0\}}$. Intuitively, we subtract from the initial rejection frequency all those parts that are due to a variation in parameters beside the one of interest. We analogously proceed for within-sample dependence.

20% ties within the sample, nominal critical p-values should be lowered by $[0.03, 0.05]^2$. The test may thus still respect a weak significance level (remember the baseline reference here being chosen $\alpha = 0.05$ ³). Maximal errors in indicated size are of 0.4 $\{g = 0, \lambda = 0.8, n = 10\}$.

Tied observations entail a loss in size-adjusted *power*, this decrease being more pronounced for larger effects. Interestingly, power declines more pronouncedly with increasing sample size for small and medium effects, though in absolute terms, the number of untied observations increases and one therefore would likely expect the reverse. Only for large effects, increasing sample size may actually help reducing power losses. Maximal power losses are of 0.54 $\{g = 0.5, \lambda = 0.8, n = 40\}$.

Within-sample dependence, our second assumption violation investigated, equally leads to an overrejection which turns out to be less pronounced than in the presence of ties only. Interestingly, the problem rather worsens with increasing sample size. For a low level of within-sample dependence $\rho_x = 0.2$, the test may still be applied by lowering critical nominal p-values by $[0.01, 0.02]$. Intermediary levels of within-sample dependence $\rho_x = 0.5$ demand a lowering of critical nominal p-values by $[0.04, 0.06]$, thus only allowing for weak significance levels (recall the nominal significance level of $\alpha = 0.05$). Maximal errors in indicated size of 0.2 occur for $\{g = 0, \rho_x = 0.8, n = 40\}$.

Power losses entailed are roughly half as important as for tied observations, its dependence on sample size displaying the same pattern. Interestingly, weak degrees of within sample dependence may actually help structuring the sample and thus improve power properties, if the effect itself is small. Maximal power losses of 0.31 occur for $\{g = 0.5, \rho_x = 0.8, n = 40\}$.

When ***both within-sample dependence and ties*** are present a mutual reinforcement of the adverse *size* effects from both violations of the test assumptions can be observed. This tendency towards overrejection declines only very slowly with increasing sample size and recognizably only for large fractions of ties and high degrees of within-sample-dependence.

² The interval is defined by the smallest and largest correction necessary for different sample sizes.

³ For an induced significance level of 0.01, the impact of assumption violations on *size* is roughly half as important as for 0.05. The corrections suggested here may thus serve as an upper bound. The impact of assumption violations on power under 0.01, however, is more severe for small and medium effects. Note that effect size and significance level are related: small and medium but yet highly significant effects alone are already difficult to identify for our small sample sizes. For the same degree of assumption violations, such effects are therefore more difficult to identify than effects of intermediate significance. Consistently, this difference in impact of assumption violations under 0.01 and 0.05 disappears for large effect sizes.

Table 1
Baseline Simulation Results

	$g=0$	$g=0$	$g=0$	$g=0.2$	$g=0.2$	$g=0.2$	$g=0.5$	$g=0.5$	$g=0.5$	$g=0.8$	$g=0.8$	$g=0.8$
	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$
$\rho_x=0$	0.051	0.050	0.050	0.084	0.129	0.183	0.232	0.306	0.554	0.749	0.857	0.990
$\lambda=0$												
$\rho_x=0$	0.085	0.088	0.086	0.091	0.127	0.206	0.244	0.316	0.503	0.666	0.773	0.986
$\lambda=0.2$												
$\rho_x=0$	0.170	0.165	0.167	0.164	0.202	0.253	0.269	0.348	0.476	0.586	0.673	0.930
$\lambda=0.5$												
$\rho_x=0$	0.460	0.359	0.337	0.336	0.456	0.374	0.385	0.527	0.516	0.565	0.610	0.837
$\lambda=0.8$												
$\rho_x=0.2$	0.058	0.061	0.057	0.064	0.100	0.193	0.246	0.313	0.564	0.744	0.855	0.998
$\lambda=0$												
$\rho_x=0.2$	0.100	0.101	0.104	0.110	0.135	0.219	0.250	0.331	0.517	0.658	0.769	0.984
$\lambda=0.2$												
$\rho_x=0.2$	0.197	0.180	0.181	0.178	0.222	0.269	0.288	0.352	0.480	0.576	0.672	0.928
$\lambda=0.5$												
$\rho_x=0.2$	0.487	0.379	0.364	0.350	0.493	0.406	0.402	0.560	0.531	0.579	0.614	0.828
$\lambda=0.8$												
$\rho_x=0.5$	0.092	0.097	0.101	0.104	0.131	0.238	0.273	0.349	0.557	0.709	0.818	0.994
$\lambda=0$												
$\rho_x=0.5$	0.134	0.143	0.147	0.145	0.164	0.246	0.280	0.347	0.523	0.643	0.751	0.970
$\lambda=0.2$												
$\rho_x=0.5$	0.245	0.232	0.231	0.227	0.265	0.293	0.318	0.404	0.505	0.585	0.665	0.917
$\lambda=0.5$												
$\rho_x=0.5$	0.525	0.446	0.415	0.411	0.532	0.446	0.451	0.589	0.557	0.591	0.634	0.821
$\lambda=0.8$												
$\rho_x=0.8$	0.211	0.246	0.239	0.246	0.250	0.332	0.367	0.401	0.555	0.671	0.745	0.964
$\lambda=0$												
$\rho_x=0.8$	0.267	0.282	0.290	0.285	0.300	0.353	0.383	0.435	0.542	0.633	0.699	0.925
$\lambda=0.2$												
$\rho_x=0.8$	0.375	0.375	0.368	0.383	0.395	0.429	0.436	0.487	0.550	0.614	0.653	0.857
$\lambda=0.5$												
$\rho_x=0.8$	0.618	0.576	0.565	0.556	0.634	0.582	0.584	0.671	0.647	0.672	0.680	0.807
$\lambda=0.8$												

Note: reported are rejection frequencies for the baseline case with $\rho_{xy} = 0.5$ and $\sigma_x = \sigma_y = 1$; test decisions are based on a significance level $\alpha = 0.05$.

Regarding *power*, large fractions of ties appearing together with high degrees of dependence entail the most severe power losses. Interestingly, power losses under simultaneous assumption violations are an increasing function of the sample size for small and medium effects. Only for large effects, power losses are again U-shapedly dependent on sample size. However, the larger the effect size, the more pronounced overall power losses.

Let us proceed with a second *scenario of unequal variances* with $\sigma_y = 4\sigma_x$ displayed in table 2. With no other assumption being violated, the test slightly underrejects the null hypothesis for $n = \{10,20\}$ and starts overrejecting it somewhat for $n \in \{30,40\}$ for all effects. Unequal variances entail power losses, depending U-shapedly on effect size. While for small effect sizes only slight power losses occur though increasing with sample size, large amounts of power are lost for medium effect size. For large effects losses disappear with increasing n , being negligible for $n \in \{20,30\}$. The *tie effect*⁴ on *size* is reinforced in this scenario. For levels of $\lambda \in \{0.2, 0.5, 0.8\}$ the test overrejects a null by maximally additional $\{0.014, 0.05, 0.07\}$ as compared to the baseline. Considering actual significance levels under this simultaneous assumption violation, a weak significance level might still be respected for $\lambda = 0.2$ where nominal critical p-values should systematically be lowered by 0.04 up to 0.05. Regarding *size-adjusted power*, unequal sample variances may substantially reduce power losses entailed by ties, this reduction being an increasing function of sample size except for large effects. Furthermore, the larger the fraction of ties, the higher the power loss reduction entailed by differences in variance.

Turning to *within-sample dependence*, its impact is weakened by unequal variances. Overrejection diminishes, entailing a need for lowering nominal critical p-values by maximally but $\{0.007, 0.014, 0.064\}$ for $\rho_x \in \{0.2, 0.5, 0.8\}$. Regarding *size-adjusted power*

⁴ The impact of a particular scenario is quantified by including one further condition in our calculations. To analyze the effect of unequal standard deviations on the impact ties show, this requires the following modification:

$$\Delta_{\lambda}^{\sigma_y=4\sigma_x} \Big|_{\{\lambda=\bar{\lambda}, n=\bar{n}\}} = ((a|\{\sigma_y=4\sigma_x\} - (a|\{\sigma_y=4\sigma_x\} - b|\{\sigma_y=4\sigma_x\})) - (c|\{\sigma_y=4\sigma_x\} - (c|\{\sigma_y=4\sigma_x\} - d|\{\sigma_y=4\sigma_x\}))) - ((a - (a - b)) - (c - (c - d)))$$

for the example above. Changes in the impact ties exert on size-adjusted power, are now calculated as follows:

$$\nabla_{\lambda}^{\sigma_y=4\sigma_x} \Big|_{\{\lambda=\bar{\lambda}, n=\bar{n}\}} = (f_r|\{\lambda=\bar{\lambda}, n=\bar{n}, \rho_x=0, \sigma_y=4\sigma_x\} - f_r|\{\lambda=0, n=\bar{n}, \rho_x=0, \sigma_y=4\sigma_x\}) - (f_r|\{\lambda=\bar{\lambda}, n=\bar{n}, \rho_x=0\} - f_r|\{\lambda=0, n=\bar{n}, \rho_x=0\}),$$

the subtrahend in either equation corresponding to the baseline case with $\sigma_x = \sigma_y$. This latter condition is not mentioned again for reasons of space.

Table 2
Results with Unequal Variances

	$g=0$	$g=0$	$g=0$	$g=0.2$	$g=0.2$	$g=0.2$	$g=0.5$	$g=0.5$	$g=0.5$	$g=0.8$	$g=0.8$	$g=0.8$
	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$
$\rho_x = 0$	0.046	0.049	0.051	0.050	0.074	0.103	0.138	0.174	0.214	0.403	0.563	0.693
$\lambda = 0$												
$\rho_x = 0$	0.091	0.100	0.096	0.098	0.126	0.157	0.179	0.206	0.263	0.404	0.536	0.640
$\lambda = 0.2$												
$\rho_x = 0$	0.215	0.201	0.197	0.200	0.237	0.234	0.271	0.282	0.352	0.432	0.523	0.599
$\lambda = 0.5$												
$\rho_x = 0$	0.525	0.426	0.400	0.398	0.542	0.436	0.431	0.443	0.575	0.531	0.563	0.593
$\lambda = 0.8$												
$\rho_x = 0.2$	0.053	0.052	0.056	0.052	0.084	0.109	0.141	0.182	0.214	0.395	0.563	0.696
$\lambda = 0$												
$\rho_x = 0.2$	0.099	0.097	0.099	0.102	0.131	0.155	0.182	0.215	0.267	0.409	0.542	0.647
$\lambda = 0.2$												
$\rho_x = 0.2$	0.227	0.209	0.205	0.206	0.252	0.255	0.263	0.288	0.350	0.430	0.516	0.603
$\lambda = 0.5$												
$\rho_x = 0.2$	0.536	0.428	0.408	0.411	0.544	0.453	0.443	0.447	0.593	0.537	0.568	0.602
$\lambda = 0.8$												
$\rho_x = 0.5$	0.059	0.063	0.065	0.068	0.087	0.121	0.152	0.183	0.230	0.409	0.554	0.682
$\lambda = 0$												
$\rho_x = 0.5$	0.114	0.120	0.115	0.117	0.144	0.171	0.200	0.220	0.270	0.418	0.539	0.644
$\lambda = 0.2$												
$\rho_x = 0.5$	0.247	0.228	0.232	0.217	0.253	0.270	0.287	0.301	0.371	0.440	0.537	0.599
$\lambda = 0.5$												
$\rho_x = 0.5$	0.558	0.458	0.440	0.435	0.565	0.477	0.462	0.461	0.607	0.558	0.569	0.612
$\lambda = 0.8$												
$\rho_x = 0.8$	0.100	0.108	0.115	0.112	0.126	0.160	0.197	0.218	0.265	0.413	0.557	0.669
$\lambda = 0$												
$\rho_x = 0.8$	0.165	0.170	0.178	0.181	0.183	0.214	0.245	0.267	0.314	0.433	0.549	0.628
$\lambda = 0.2$												
$\rho_x = 0.8$	0.294	0.287	0.293	0.297	0.309	0.324	0.338	0.334	0.405	0.471	0.531	0.603
$\lambda = 0.5$												
$\rho_x = 0.8$	0.586	0.513	0.494	0.496	0.592	0.519	0.517	0.514	0.628	0.588	0.606	0.632
$\lambda = 0.8$												

Note: reported are rejection frequencies for the baseline case with $\rho_{xy} = 0.5$, $\sigma_x = 1$ and $\sigma_y = 4$; test decisions are based on a significance level $\alpha = 0.05$.

slight improvements are observed for $\rho_x \in \{0.2, 0.5\}$ while under $\rho_x = 0.8$, power gains as compared to the scenario with equal variances reach up to 0.23. Turning to a last constellation where ***both ties and within-sample dependence*** are present, the tie impact on *size* under unequal variances dominates for small levels of within sample dependence. That is, overrejection increases as compared to the baseline.

However, for accruing within-sample dependence the latter effect takes lead and overrejection for simultaneous assumption violations under unequal sample variances diminishes in comparison to the baseline. Regarding the impact on *power*, an overall improvement is observed. Thus, power losses decrease as compared to the reference case and more pronouncedly do so for both increasing within-sample dependence and tie fraction. However, power losses decrease decisively faster with increasing within-sample dependence.

A ***third scenario*** as shown in table 3 allows us to investigate a possible impact of the very ***nature of within-sample dependence*** and for this purpose introduces an MA rather than an AR structure. While for $\rho_x = 0.2$, the impact on *size* is manifested by a slightly increased overrejection, the reverse holds for stronger within-sample dependence: in presence of an MA structure overrejection diminishes by $[0.02; 0.03]$ for $\rho_x = 0.5$ as compared to an AR structure. *Power* improves for the same degree of within-sample dependence up to 0.046 and increasingly does so for larger sample sizes. Again, power improvements become more visible for $\rho_x = 0.5$ and display an increasing dependence on sample size. Thus, the kind of within-sample dependence equally affects the relationship between power and sample size. Furthermore, when ***both ties and within-sample dependence*** are present, power properties improve slightly with increasing g . riven by this effect, overrejection entailed by simultaneous assumption violation slightly declines.

A ***fourth scenario*** associated with table 4 subsequently investigates the impact of the ***nature of ties*** by introducing one large tie instead of many small. With respect to *size*, the impact is tremendous. For the smallest fraction $\lambda = 0.2$, overrejection entailed by ties increases by additional $[0.03, 0.05, 0.08, 0.07]$ for $n \in \{10, 20, 30, 40\}$. Thus, weak significance levels cannot even be respected for the smallest fraction of ties investigated. Rejection frequency with samples displaying no difference reaches certainty for $\lambda = 0.8$. In presence of ***both ties and within-sample dependence***, overrejection increases beyond the tie-driven effect by $[0.04, 0.07, 0.08, 0.09]$ for $n \in \{10, 20, 30, 40\}$.

Table 3
Results for Within-Sample Dependence as MA(1) Instead of AR(1)

	$g=0$	$g=0$	$g=0$	$g=0.2$	$g=0.2$	$g=0.2$	$g=0.5$	$g=0.5$	$g=0.5$	$g=0.8$	$g=0.8$	$g=0.8$
	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$	$n=10$	$n=20$	$n=30$
$\rho_x = 0$	0.052	0.049	0.048	0.054	0.093	0.131	0.190	0.227	0.302	0.554	0.746	0.858
$\lambda = 0$												
$\rho_x = 0$	0.091	0.083	0.087	0.088	0.123	0.164	0.199	0.247	0.308	0.513	0.665	0.781
$\lambda = 0.2$												
$\rho_x = 0$	0.170	0.166	0.173	0.168	0.204	0.228	0.255	0.285	0.350	0.461	0.584	0.669
$\lambda = 0.5$												
$\rho_x = 0$	0.451	0.358	0.343	0.328	0.462	0.380	0.378	0.384	0.528	0.514	0.560	0.606
$\lambda = 0.8$												
$\rho_x = 0.2$	0.062	0.059	0.059	0.057	0.097	0.147	0.190	0.240	0.309	0.564	0.737	0.859
$\lambda = 0$												
$\rho_x = 0.2$	0.092	0.096	0.102	0.095	0.140	0.181	0.214	0.255	0.324	0.510	0.663	0.766
$\lambda = 0.2$												
$\rho_x = 0.2$	0.186	0.174	0.179	0.176	0.218	0.232	0.255	0.294	0.359	0.470	0.580	0.660
$\lambda = 0.5$												
$\rho_x = 0.2$	0.474	0.382	0.352	0.352	0.486	0.406	0.403	0.397	0.557	0.535	0.576	0.615
$\lambda = 0.8$												
$\rho_x = 0.5$	0.071	0.073	0.078	0.081	0.114	0.158	0.203	0.256	0.326	0.568	0.730	0.834
$\lambda = 0$												
$\rho_x = 0.5$	0.117	0.110	0.118	0.114	0.150	0.192	0.233	0.270	0.335	0.510	0.657	0.760
$\lambda = 0.2$												
$\rho_x = 0.5$	0.210	0.201	0.195	0.199	0.244	0.256	0.272	0.305	0.374	0.480	0.587	0.660
$\lambda = 0.5$												
$\rho_x = 0.5$	0.504	0.407	0.395	0.373	0.524	0.430	0.421	0.414	0.562	0.541	0.576	0.618
$\lambda = 0.8$												

Note: reported are rejection frequencies for the baseline case with $\rho_{xy} = 0.5$, $\sigma_x = 1$ and $\sigma_y = 1$; test decisions are based on a significance level $\alpha = 0.05$.

Table 4
Results for the Variant with One Big (Median) Tie Instead of Many Small Ties

	$g=0$	$g=0$	$g=0$	$g=0.2$	$g=0.2$	$g=0.2$	$g=0.5$	$g=0.5$	$g=0.5$	$g=0.8$	$g=0.8$	$g=0.8$
	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$
$\rho_x = 0$	0.054	0.048	0.047	0.048	0.086	0.132	0.180	0.231	0.301	0.558	0.748	0.861
$\lambda = 0$	0.054	0.048	0.047	0.048	0.086	0.132	0.180	0.231	0.301	0.558	0.748	0.861
$\rho_x = 0$	0.119	0.135	0.159	0.161	0.171	0.251	0.314	0.375	0.411	0.664	0.812	0.892
$\lambda = 0.2$	0.119	0.135	0.159	0.161	0.171	0.251	0.314	0.375	0.411	0.664	0.812	0.892
$\rho_x = 0$	0.355	0.512	0.591	0.668	0.391	0.573	0.668	0.743	0.559	0.801	0.878	0.940
$\lambda = 0.5$	0.355	0.512	0.591	0.668	0.391	0.573	0.668	0.743	0.559	0.801	0.878	0.940
$\rho_x = 0$	0.714	0.968	0.999	1.000	0.720	0.969	1.000	1.000	0.769	0.978	0.999	1.000
$\lambda = 0.8$	0.714	0.968	0.999	1.000	0.720	0.969	1.000	1.000	0.769	0.978	0.999	1.000
$\rho_x = 0.2$	0.057	0.060	0.065	0.062	0.098	0.148	0.192	0.247	0.318	0.551	0.732	0.846
$\lambda = 0$	0.057	0.060	0.065	0.062	0.098	0.148	0.192	0.247	0.318	0.551	0.732	0.846
$\rho_x = 0.2$	0.146	0.165	0.182	0.199	0.191	0.269	0.332	0.391	0.412	0.659	0.806	0.891
$\lambda = 0.2$	0.146	0.165	0.182	0.199	0.191	0.269	0.332	0.391	0.412	0.659	0.806	0.891
$\rho_x = 0.2$	0.373	0.519	0.597	0.668	0.407	0.592	0.671	0.753	0.572	0.791	0.876	0.930
$\lambda = 0.5$	0.373	0.519	0.597	0.668	0.407	0.592	0.671	0.753	0.572	0.791	0.876	0.930
$\rho_x = 0.2$	0.713	0.968	0.999	1.000	0.729	0.968	1.000	1.000	0.775	0.978	0.999	1.000
$\lambda = 0.8$	0.713	0.968	0.999	1.000	0.729	0.968	1.000	1.000	0.775	0.978	0.999	1.000
$\rho_x = 0.5$	0.093	0.095	0.106	0.101	0.145	0.180	0.226	0.280	0.339	0.554	0.719	0.820
$\lambda = 0$	0.093	0.095	0.106	0.101	0.145	0.180	0.226	0.280	0.339	0.554	0.719	0.820
$\rho_x = 0.5$	0.185	0.217	0.233	0.245	0.227	0.302	0.367	0.417	0.435	0.640	0.780	0.855
$\lambda = 0.2$	0.185	0.217	0.233	0.245	0.227	0.302	0.367	0.417	0.435	0.640	0.780	0.855
$\rho_x = 0.5$	0.424	0.569	0.646	0.715	0.456	0.620	0.699	0.767	0.595	0.785	0.868	0.921
$\lambda = 0.5$	0.424	0.569	0.646	0.715	0.456	0.620	0.699	0.767	0.595	0.785	0.868	0.921
$\rho_x = 0.5$	0.734	0.972	0.999	1.000	0.748	0.975	1.000	1.000	0.792	0.979	1.000	1.000
$\lambda = 0.8$	0.734	0.972	0.999	1.000	0.748	0.975	1.000	1.000	0.792	0.979	1.000	1.000
$\rho_x = 0.8$	0.209	0.235	0.242	0.242	0.247	0.301	0.328	0.357	0.408	0.560	0.671	0.748
$\lambda = 0$	0.209	0.235	0.242	0.242	0.247	0.301	0.328	0.357	0.408	0.560	0.671	0.748
$\rho_x = 0.8$	0.336	0.371	0.387	0.399	0.362	0.419	0.467	0.495	0.492	0.631	0.724	0.796
$\lambda = 0.2$	0.336	0.371	0.387	0.399	0.362	0.419	0.467	0.495	0.492	0.631	0.724	0.796
$\rho_x = 0.8$	0.540	0.684	0.736	0.777	0.557	0.703	0.759	0.807	0.634	0.790	0.851	0.896
$\lambda = 0.5$	0.540	0.684	0.736	0.777	0.557	0.703	0.759	0.807	0.634	0.790	0.851	0.896
$\rho_x = 0.8$	0.784	0.983	1.000	1.000	0.792	0.986	0.999	1.000	0.808	0.988	1.000	1.000
$\lambda = 0.8$	0.784	0.983	1.000	1.000	0.792	0.986	0.999	1.000	0.808	0.988	1.000	1.000

Note: reported are rejection frequencies for the baseline case with $\rho_{xy} = 0.5$, $\sigma_x = 1$ and $\sigma_y = 1$; test decisions are based on a significance level $\alpha = 0.05$.

The impact on *power* is less unanimous. For weak effects of $g = 0.2$, a small fraction $\lambda = 0.2$ of a larger tie helps structuring the sample, and does so more efficiently than smaller ties improving power by additional $[0.01, 0.03, 0.04, 0.06]$ for $n \in \{10, 20, 30, 40\}$. Power gains for $g = 0.2$ as compared to small ties are furthermore increasing with sample size. However, larger fractions of large ties entail additional power losses for all situations both with increasing sample size and drastically so for larger effect sizes. Maximally incurred additional power losses of 0.5 are observed for $\{g = 0.8, n = 40, \lambda = 0.8\}$. **Both ties and within-sample dependence** yield subadditively increased power losses for low fraction of ties, for high fraction of ties however, power losses exceed the sum of the two single effects.

A **fifth scenario** displayed in table 5 sheds light on the impact of ties and within-sample dependence under asymmetry of the underlying distributions. In the absence of any further assumption violation and with the samples showing no difference, asymmetry entails a tendency towards underrejection. Its effect on *size* is of approximately -0.02. *Power* under asymmetry increases overall, the improvement being dependent on both effect and sample size. Maximal power gains of 0.07 occur for $\{g = 0.5, n = 30\}$. The **tie effect** on *size* is ambiguously influenced by asymmetry. Generally, entailed overrejection decreases, and especially does so for the case of interest with $\lambda = 0.2$ where the effect is slightly weakened but for $n = 10$. Regarding the tie effect on *power*, asymmetry yields power improvements of up to 0.07 $\{g = 0.8, n = 40, \lambda = 0.8\}$. A more pronounced change is observed for the impact of **within-sample dependence**. Its effect on *size* is overcompensated for $\rho_x = 0.2$ and $\rho_x = 0.5$. Depending on sample size, it diminishes by $[0.004, 0.005, 0.006, 0.008]$ for the former and by $[0.04, 0.05, 0.05, 0.06]$ for the latter case. Furthermore, asymmetry causes ρ_x – entailed overrejection to decline much faster with increasing sample size. Even for $\rho_x = 0.8$, starting with an actual size of 0.12, the test reaches nominal size for $n = 40$. *Power* losses entailed by within-sample dependence decrease quite substantially under asymmetry, and more pronouncedly do so with increasing sample size. Maximal power gains of 0.21 occur for $\{g = 0.8, \rho_x = 0.8, n = 40\}$. Only in the absence of any effect and for very small samples $n \in \{10, 20\}$ slight *power* losses are observed in comparison to the baseline.

Table 5
Results for Gamma-Distributed Instead of Normal-Distributed Series

	$g=0$	$g=0$	$g=0.2$	$g=0.2$	$g=0.2$	$g=0.5$	$g=0.5$	$g=0.5$	$g=0.8$	$g=0.8$	$g=0.8$					
	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$	$n=10$	$n=20$	$n=30$	$n=40$				
$\rho_x = 0$	0.029	0.028	0.029	0.030	0.065	0.114	0.170	0.234	0.289	0.592	0.803	0.915	0.641	0.952	0.997	0.999
$\lambda = 0$																
$\rho_x = 0$	0.064	0.065	0.065	0.069	0.098	0.157	0.212	0.271	0.312	0.543	0.729	0.843	0.613	0.884	0.977	0.994
$\lambda = 0.2$																
$\rho_x = 0$	0.153	0.138	0.147	0.139	0.191	0.219	0.265	0.311	0.358	0.506	0.642	0.736	0.579	0.792	0.909	0.960
$\lambda = 0.5$																
$\rho_x = 0$	0.435	0.338	0.311	0.299	0.454	0.367	0.369	0.395	0.528	0.531	0.595	0.659	0.652	0.717	0.819	0.884
$\lambda = 0.8$																
$\rho_x = 0.2$	0.031	0.029	0.030	0.028	0.067	0.116	0.175	0.237	0.284	0.596	0.796	0.911	0.650	0.953	0.996	1.000
$\lambda = 0$																
$\rho_x = 0.2$	0.064	0.062	0.061	0.058	0.104	0.150	0.212	0.263	0.316	0.549	0.722	0.836	0.618	0.884	0.972	0.993
$\lambda = 0.2$																
$\rho_x = 0.2$	0.149	0.139	0.142	0.138	0.191	0.227	0.262	0.302	0.344	0.513	0.636	0.740	0.576	0.789	0.907	0.965
$\lambda = 0.5$																
$\rho_x = 0.2$	0.451	0.332	0.314	0.299	0.457	0.378	0.369	0.385	0.548	0.541	0.603	0.653	0.659	0.724	0.811	0.880
$\lambda = 0.8$																
$\rho_x = 0.5$	0.032	0.029	0.027	0.026	0.068	0.108	0.164	0.219	0.295	0.571	0.775	0.890	0.626	0.939	0.995	1.000
$\lambda = 0$																
$\rho_x = 0.5$	0.070	0.063	0.059	0.061	0.102	0.151	0.195	0.243	0.317	0.522	0.693	0.813	0.597	0.861	0.963	0.992
$\lambda = 0.2$																
$\rho_x = 0.5$	0.166	0.142	0.137	0.131	0.199	0.216	0.248	0.280	0.379	0.493	0.602	0.703	0.586	0.773	0.893	0.948
$\lambda = 0.5$																
$\rho_x = 0.5$	0.482	0.356	0.317	0.307	0.504	0.393	0.389	0.377	0.564	0.547	0.593	0.644	0.668	0.729	0.804	0.860
$\lambda = 0.8$																
$\rho_x = 0.8$	0.119	0.074	0.059	0.048	0.157	0.162	0.169	0.202	0.342	0.508	0.660	0.784	0.586	0.837	0.955	0.991
$\lambda = 0$																
$\rho_x = 0.8$	0.176	0.131	0.113	0.091	0.215	0.199	0.214	0.233	0.357	0.476	0.586	0.687	0.576	0.767	0.890	0.954
$\lambda = 0.2$																
$\rho_x = 0.8$	0.303	0.257	0.218	0.194	0.335	0.311	0.291	0.284	0.431	0.479	0.535	0.604	0.584	0.705	0.799	0.865
$\lambda = 0.5$																
$\rho_x = 0.8$	0.599	0.523	0.484	0.462	0.596	0.531	0.509	0.490	0.641	0.614	0.613	0.633	0.708	0.711	0.751	0.795
$\lambda = 0.8$																

Note: reported are rejection frequencies for the baseline case with $\rho_{xy} = 0.5$, $\sigma_x = \sigma_y = 1$ and skewness of one; test decisions are based on a significance level $\alpha = 0.05$.

