

Sequential Specification Tests to Choose a Model: A Change-Point Approach

Adam Sales

January 27, 2017

1 Introduction

One of many mis-matches between best practices recommended by statisticians and practice in quantitative research regards model selection. Statisticians conceptualize model selection as a tradeoff between bias and variance. Many quantitative researchers think about model selection as choosing the best model that satisfies the assumptions of their intended statistical test or estimator—essentially minimizing variance while constraining bias at zero. This latter outlook leads researchers towards hypothesis tests of model assumptions; in particular, a sequence of hypothesis tests, for a sequence of models, ordered by preferability. The best model whose assumptions “pass” a hypothesis test is chosen.

Do hypothesis tests make any sense in model selection? For one, “all

models are wrong” (? , p. 2) and “there’s no such thing as unbiased estimation” (?), so the search for a correct model might be hopeless, and therefore pointless. Further, the logic of null-hypothesis testing seems incompatible with this framework. The results of a null hypothesis test, of course, are never evidence in favor of a null hypothesis—null hypotheses can only be rejected, not accepted.

On the other hand, “some models are useful,” and depending on their intended use, their usefulness may depend on *approximate* correctness. If so, hypothesis tests may have a role to play. Specification tests already exist for most common models, and they are regularly taught in introductory quantitative methods classes. If their use in model selection could be made conceptually sound, they are likely to be actually used—and maybe even correctly.

This paper will borrow a clever idea from change-point or threshold estimation to the more general problem of model selection from hypothesis tests—more accurately, p-values. ? points out that in a process with a change point, the p-values from a sequence of tests of a null regression function are uniformly-distributed as long as the regression function is correct, but asymptotically zero when the function is not correct. They use this dichotomous behavior to construct a simple, consistent estimator of the change-point—the point at which the null model stops being correct.

In the same way, their estimator can choose the change-point in a sequence of models, when models stop being correct. In doing so, it shifts

the model selection rationale away from the logic of hypothesis testing and towards the logic of estimation. In the tradition of constructing confidence intervals from hypothesis tests and [??](#), their estimator exploits the behavior of hypothesis tests to estimate quantities of interest. Further, as opposed to model selectors based on strict hypothesis-testing logic, an individual test result will itself not drive the change-point estimator, which is instead based on the entire sequence of p-values. Thus, the change-point view of model selection is arguably conceptually more satisfying and practically more reliable than the conventional test-based approach.

The following sub-section will briefly introduce two running examples of sequential specification tests: choosing a bandwidth for a regression discontinuity design and choosing a lag order for a time-series model. Next, [section ??](#) will review the formalism of SSTs and discuss common SST-based model selectors. [Section 3](#) will introduce the new method, [section 4](#) will demonstrate some of its properties in a simulation study, [5](#) will apply it to the running examples, and [6](#) will conclude.

1.1 SSTs in Regression Discontinuity and Time Series

[Figure ??](#) displays two datasets that will serve as illustrations of SSTs. [Section 5](#) will discuss both of these examples in more detail. The brief overview here will be helpful to fix ideas.

[Figure ??A](#) plots data that [?](#) (LSO) used to estimate the effect of academic probation. Students at an unnamed large Canadian University were

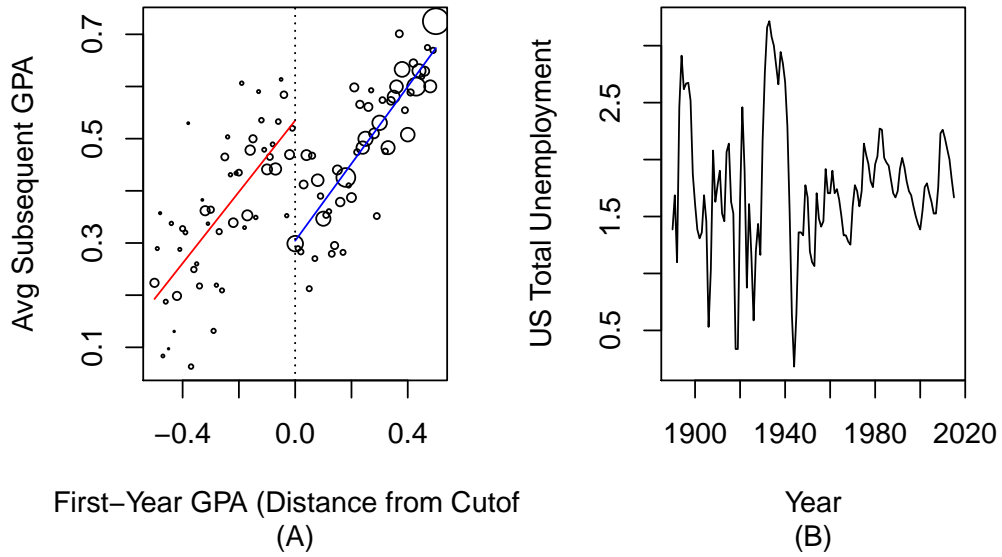


Figure 1: Two data examples for SSTs. Plot (A) shows data from ?—subsequent grade point averages (GPAs) for students at a large Canadian university, as a function of first-year GPAs. Subsequent GPAs are averaged by first-year GPAs, which are centered at the academic probation cutoff (dotted line), and the sizes of the plotted points are proportional to the number of students with each first-year GPA. The red and blue lines are linear least-squares fits on either side of the cutoff. Students with first-year GPAs to the left of the cutoff are put on probation. Plot (B) shows a time-series of log annual United States total unemployment from 1890 to 2015. Data were combined from ? and ?.

put on academic probation—simultaneously given extra help and threatened with suspension—if their first-year cumulative grade point averages (GPAs) fell below a cutoff. This is an example of a regression discontinuity design (RDD) (??), in which treatment (in this case academic probation) is assigned if a numeric “running variable” R (first-year GPA) falls below (or above) a pre-specified cutoff c . Typically (e.g. ???) analysts will fit regression models $Y = f_1(R) + \epsilon$ and $Y = f_2(R) + \epsilon$ to data on either side of the cutoff, modeling the relationship between R and an outcome of interest Y . The difference between the models’ predictions when R is set equal to the c is interpreted as a “local average treatment effect,” roughly speaking the treatment effect when the running variable is equal to the cutoff (?). Figure ??A shows one of the outcomes LSO considered, students’ subsequent GPAs, along with linear regression models below (in red) and above (blue) c , which is signified with a dotted line. A simpler alternative approach, suggested in ?, models the relationships between Y and R on either side of c as constant, and treats the data as if they were generated by a randomized experiment.

Of course, misspecified regression models will lead to biased treatment effect estimators. To minimize the influence of model misspecification, researchers will typically fit the regression models using only subjects for whom $R \in \mathcal{W}_b \equiv (c - b, c + b)$ for some bandwidth $b > 0$. A number of options exist for choosing the RDD bandwidth, including cross-validation (?) and asymptotic minimization of mean-squared-error (?). ? and others ?? suggest SSTs of covariate balance—at a sequence of candidate bandwidths b ,

test for the presence of a “treatment effect” on a pre-treatment covariate X , referred to as covariate imbalance. A window choice \mathcal{W}_b that, paradoxically, leads to a statistically significant non-zero treatment effect on a covariate is unacceptable; on the other hand, larger windows include larger datasamples, yielding higher precision. Therefore, SSTs could be used to choose the largest b for which a hypothesis test fails to reject the hypothesis of covariate balance.

Figure ??B shows the annual total unemployment rate in the United States from 1890 to 2015. One of the simpler models for time series such as these is an order p autoregression, or $AR(p)$ under which the value of the time series at point t may depend on its historical values at $t - 1, \dots, t - p$ but, conditional on those, is independent of values at points before $t - p$. SSTs can be useful here, too: researchers may test model fit for a sequence of lag orders p , and choose the smallest p that the tests fail to reject. Here a smaller lag orders p are preferable since they lead to more parsimonious models and more precise estimates.

2 The Setup, in General

Say, in specifying a model, a research must choose from a discrete, ordered, set of specifications $d = 1, 2, \dots, D$. The resulting model must satisfy testible assumption \mathcal{A} . Assume that either \mathcal{A} is false for all d , or that for some $1 \leq d^* \leq D$, \mathcal{A} is true for $d \leq d^*$ and false for all $d > d^*$. Further

assume that if d^* exists, it is the optimal choice—for instance, the smallest model, or the biggest dataset, that satisfies \mathcal{A} . Finally, assume the researcher has chosen a valid, unbiased test of \mathcal{A} and calculated p-values for each d : $\mathbf{p}_D = p_1, \dots, p_d, \dots, p_D$. The procedure here is to use \mathbf{p}_D to choose a specification \hat{d} that is as large as possible without violating \mathcal{A} .

A common choice for d in this scenario relies on the logic of null hypothesis testing: for a pre-specified $\alpha \in (0, 1)$, let

$$\bar{d}_\alpha \equiv \max\{d : p_d > \alpha\}.$$

That is, \bar{d}_α is the largest value of d for which the null hypothesis that \mathcal{A} is true for $d \leq \bar{d}_\alpha$ cannot be rejected at level α . Although it may seem as though the multiplicity of tests involved in this procedure invalidates the null hypothesis framework, it turns out that this is not the case: the “stepwise intersection-union principal” ??? insures that the family-wise error rate is maintained. That is, the probability of falsely rejecting the null—choosing $\bar{d}_\alpha < d^*$, is bounded by α . \bar{d}_α is the specification that would result from testing null hypotheses backwards: for $d' = D, D-1, \dots, d, \dots, 1$, test $H_{0d'} : \mathcal{A}$ is true for $d \leq d'$. Then, stop testing at $d' = \bar{d}_\alpha - 1$ —the first d' for which $p_{d'} \geq \alpha$; reject all null hypotheses $H_{0d'}$ for which $d' \geq \bar{d}_\alpha$, and fail to reject the rest. This protects the family-wise error rate of α since rejecting *any* true null implies rejecting the first true null—a probability α event.

Another common choice for \hat{d} , say $\underline{\hat{d}}_\alpha$, does not have this property. Let

$$\underline{\hat{d}}_\alpha \equiv \min\{d : p_d < \alpha\} - 1 \quad (1)$$

$\underline{\hat{d}}_\alpha$ selects \hat{d} to be the largest value of d before the first significant p-value. This is equivalent to the opposite procedure as $\bar{\hat{d}}_\alpha$: start with the $d' = 1$ and test sequentially for larger values of d' until the first rejection, at $\underline{\hat{d}}_\alpha$, then stop; reject all null hypotheses $H_{0d'}$ for $d' \geq \underline{\hat{d}}_\alpha$ and fail to reject the rest. This procedure does not control family-wise error rates—it is likely to reject more than $100\alpha\%$ valid specifications.

This paper will focus on two data scenarios for SSTs, corresponding to the two examples in Section ???. In the first, the SSTs help determine which data are included in the analysis. For instance, choosing the bandwidth of a regression discontinuity design, or choosing the parameters of a matching design. In this scenario, each choice d corresponds to rows in the dataset that could be included in the analysis. Formally, let $\mathcal{I} = \{1, \dots, N\}$, indices for N candidate cases to be fit in a model. Then let $\mathcal{I} = \rangle_1 \cup \rangle_2 \cup \dots \cup \rangle_d \cup \dots \text{union} \rangle_D$. The choice of \hat{d} means fitting the model to the dataset including each of these subsets, $\mathcal{I}_d = \rangle_1 \cup \dots \rangle_d$. Note that the sets denoted with lower-case \rangle_d are disjoint, $\rangle_d \cap \rangle_{d'} = \emptyset$, those denoted with upper-case \mathcal{I}_d are nested— $d > d'$ implies $\mathcal{I}_{d'} \subset \mathcal{I}_d$, and the full set of indices, noted without a subscript, $\mathcal{I} = \mathcal{I}_D$. Finally, let $n_d = |\rangle_d|$, where $|\cdot|$ denotes cardinality, and $\bar{n}_d = |\mathcal{I}_d|$.

A second scenario applies when the dataset is fixed, but the model is

not. Here, d indexes a *pre-specified* sequence of models. For instance, using SSTs to choose the lag order p in an $AR(p)$ time series model. Then let \mathbf{X} denote the full set of variables, \mathbf{x}_d denote the set of variables that would be *subtracted* in the d^{th} step of the sequence, and $\mathbf{X}_d = \mathbf{x}_{d+1} \cup \dots \cup \mathbf{x}_D$ denote the set of variables that would be included in the analysis, were the analyst to choose d . Note here that bigger values of d correspond to smaller models. In this scenario, the sample size is fixed at N .

2.1 Model Selection and the Logic of Null Hypothesis Testing

In order to avoid certain methodological mistakes, it may be helpful to clarify some of the conceptual distinctions between SSTs and conventional null hypothesis tests (NHTs). The logic of NHTs is familiar to anyone who has taken (and understood) even the most basic college statistics course; nonetheless we restate it here to distinguish it from the logic of SSTs. Typically, researchers use NHTs to reject a null hypothesis that they consider uninteresting—most of the time, that a model parameter is equal to zero—and interpret rejection as evidence in favor of an interesting alternative hypothesis. NHTs cap the probability of a type-I error—falsely rejecting a true null hypothesis—and, given that constraint, seek to minimize the probability of a type-II error, failing to reject a false null hypothesis.

SSTs reverse some of these elements; most importantly, the goal of SSTs

is to identify specifications in which an assumption \mathcal{A} is plausible, rather than to identify true alternative hypothesis. In the same vein, type-II errors are typically of more concern for SSTs than for typical NHTs, and type-I errors are less problematic. In fact, a type-II error from a specification test could lead a researcher to fit a misspecified model, which in turn may inflate the probability of a type-I error in her final outcome analysis. For that reason, some methodologists recommend setting α substantially higher for specification tests than for NHTs in outcome analyses. Still, the hypothesis testing framework, in the case of point null hypotheses, does not allow a researcher to fix the type-II error rate at a pre-specified value, and then optimize the type-I error rate, though that might be ideal for specification tests.

In fact, in continuous data models with continuous parameter spaces, no hypothesis test can provide any evidence in favor of a point null hypothesis. For instance, take the common $H_0 : \theta = 0$, for some parameter $\theta \in \mathbb{R}$. In finite samples, for any type-I or type-II error rate, there will always be some plausible alternative hypothesis $H_a \theta = \epsilon \neq 0$. Further, in these situations, finite sample estimates $\hat{\theta}$ will almost surely be non-zero. This is important to state to avoid misinterpretations of SST procedures as providing evidence, or showing, that an assumption \mathcal{A} is true for certain specifications d . A common Bayesian argument (e.g. ?, p. 439; ?) states that, theoretically, nearly all null hypotheses are false anyway—so testing them makes little sense. In the case of specification tests, that means that an assumption \mathcal{A} can be assumed

to be false for all d without even conducting a test; in other words, “all models are wrong” (? , p. 2).

“But some are useful.” In practice there is much to be gained by considering assumptions such as \mathcal{A} . In this framework, it may indeed make sense to identify a set of specifications d for which \mathcal{A} is plausible, or approximately true, and SSTs can be useful in this regard—as long as they are understood correctly, and not as providing evidence *for* \mathcal{A} .

In many scenarios the choice of d involves a bias-variance tradeoff: if $d > d^*$, then \mathcal{A} is false and the resulting analysis will be biased. On the other hand, a sub-optimal choice for d often means a high-variance estimate. For instance, in the RDD bandwidth case, choosing $d > d^*$ might mean fitting a misspecified model to Y and R , but choosing $d \ll d^*$ means discarding data that can boost precision. Rather than choosing a criterion, such as mean-squared-error, that balances bias and variance, the SST approach may be seen as an attempt to hold bias at approximately zero, and minimize variance under that constraint. Granted, this is obviously an overly-optimistic take on model fitting; still, SSTs hope to constrain bias to be approximately zero, and from there minimize variance.

2.2 More Reservations with Null Hypothesis Testing for Model Selection

Applying a strict hypothesis-testing framework to SSTs for model selection has some additional drawbacks. First, it requires researchers to choose a test-level α . While using tuning parameters to mediate the bias-variance tradeoff is not uncommon in statistics, the level α is a particularly hard parameter to choose.

? poses an additional problems with the use of hypothesis tests to choose a model: the need to specify a null hypothesis. In their words (p. 179),

Whenever a hypothesis test is used to choose between two models, one model must be selected as a null hypothesis. In most instances, this is usually the more parsimonious model and typically a nested test is applied. Often it is difficult to distinguish between the two models because of data quality (multicollinearity, near-identification, or the models being very similar such as in testing for integration). In such cases, the model chosen to be the null hypothesis is unfairly favored.

In other words, because of the structure of null hypothesis tests, which constrain the type-I error rate, the null model is unfairly favored. In our terminology, \hat{d} is likely to be too small, perhaps $\mathbb{E}\hat{d} < d^*$. However, such a bias (if it indeed exists) needn't doom SSTs—an underestimated \hat{d} is merely sub-optimal. In our setup, choosing \hat{d} to be too low will yield and ineffi-

cient, but still valid, model. Would that every statistical model were valid yet suboptimal!

More broadly, perhaps, one might argue that null hypothesis tests are design to rule out hypotheses that are inconsistent with the data, not to estimate parameters. However, as ? showed, these aims are not contradictory—tests that rule out implausible hypothesis may also point researchers towards the correct answer.

Moving from rejecting implausible specifications to estimating optimal specifications requires a theory, or at least a reasonable heuristic. The following section will suggest one.

3 Finding the Change-Point

In the context of change point estimation, ? suggests such a heuristic. They discuss a random variable x_t , whose distribution is a function of a continuous covariate t . For $t < d^*$, $\mathbb{E}x_t = \tau_0$, a constant; for $t > d^*$, $\mathbb{E}x_t > \tau_0$. They propose an estimate of d_0 based on p-values p_t testing the hypotheses $H_{0t} : \mathbb{E}x_t = \tau_0$. They note that for $t < d^*$, the null hypotheses are true, so $p_t \sim U(0, 1)$, and $\mathbb{E}p_t = 1/2$; when $t > d^*$, the null hypotheses are false, and the p-values converge in probability to zero. That fact leads them to the following least-squares estimator for d^* :

$$\hat{d}_M \equiv \arg \min_{d \in \mathbb{N}} \sum_{t \leq d} (p_t - 1/2)^2 + \sum_{t > d} p_t^2.$$

In other words, the estimate \hat{d}_M is the point at which the p-values cease behaving as p-values testing a true null, with mean $1/2$, and instead are drawn from a distribution with a lower mean. It turns out that an equivalent expression for \hat{d}_M is:

$$\hat{d}_M = \operatorname{argmax}_d \sum_{t \leq d} (p_t - 1/4). \quad (2)$$

? shows that as n_t , the number of data points at each value t , and the number of sampled values of t increase, \hat{d}_M converges in probability to d^* .

The same broad logic applies to any set of p-values from sequential tests: $\hat{d}_M = \operatorname{argmax}_d \sum_{t \leq d} (p_t - 1/4)$ may be considered an estimate of d^* . In the case of SSTs, for $d \leq d^*$, p-values p_d are draws from a $U(0, 1)$ distribution, and hence have mean $1/2$, and, as n_d or N increase, $p_d \rightarrow_p 0$ for $d > d^*$. Some differences in the details, though, lead to differences in \hat{d}_M 's behavior. For instance:

Lemma 1. *If indeed $p_d \rightarrow_p 0$ for $d > d^*$, as n_d or N increase, then \hat{d}_M is asymptotically conservative: $\Pr(\hat{d}_M > d^*) \rightarrow 0$.*

Proof. For each d , $\Pr(p_d - 1/4 > 0) \rightarrow 0$, implying that for all d' , $\Pr(\sum_{d^* < t \leq d'} (p_t - 1/4) > 0) \rightarrow 0$. Therefore, for $d^* < d \leq D$, $\Pr(\sum_{t \leq d} (p_t - 1/4) > \sum_{t \leq d^*} (p_t - 1/4)) \rightarrow 0$. \square

That is, as sample size increases, the probability that \hat{d}_M suggests a model that violates assumption \mathcal{A} decreases to zero. The same property holds for

\bar{d}_α , with $\alpha > 0$ fixed, for the same reason.

On the other hand, even with an infinite sample \hat{d}_M may choose a sub-optimal model, $\hat{d}_M < d^*$. As sample size grows, the distribution of p_d , $d \leq d^*$ remains stable at $U(0,1)$. When $p_d^* - 1/4 < 0$, $\hat{d}_M \neq d^*$, since $\sum_{d \leq d^*-1} (p_d - 1/4) > \sum_{d \leq d^*} (p_d - 1/4)$. Since $Pr(p_d^* - 1/4 < 0) = 1/4$ regardless of sample size, \hat{d}_M will be conservative in large samples. The difference between the SST case discussed here and the change-point case in ? is that the latter case relies on a continuous covariate that may be sampled from any point on the unit interval, whereas in the SST case the choice set $d = 1, 2, \dots, D$ is discrete and held fixed in the asymptotics.

In a way, \hat{d}_M is similar to $\bar{d}_{0.25}$, the largest d for which $p_d > \alpha = 0.25$, since both penalize p-values lower than 0.25. However, for a given set of p-values, $\hat{d}_M \leq \bar{d}_{0.25}$. To see this, note that for all $d > \bar{d}_{0.25}$, $p_d < 0.25$, so every summand $p_d - 1/4$ after $\bar{d}_{0.25}$ is negative. Therefore, the maximum of $\sum_{i \leq d} (p_i - 1/4)$ must occur with $d \leq \bar{d}_{0.25}$. While $\bar{d}_{0.25}$ and \hat{d}_M may often coincide, there are also cases in which $\hat{d}_M < \bar{d}_{0.25}$. This will happen when the maximum value of the random walk in (??), occurs prior to $\bar{d}_{0.25} - 1$. Then, \hat{d}_M will only equal $\bar{d}_{0.25}$ if $p_{\bar{d}_{0.25}} - 1/4 > \max_d \{ \sum_{i \leq d} (p_i - 1/4) \} - \sum_{i \leq \bar{d}_{0.25}-1} (p_i - 1/4)$.

In general, the difference between \bar{d}_α and \hat{d}_M will be most pronounced when the distributions of p-values for $d > d^*$ are not monotonically decreasing in probability—in such a scenario, it is most probable that an errent p-value for $d \gg d^*$ will be greater than α ; one p-value determines \bar{d}_α , but \hat{d}_M relies

on the entire set of p-values.

3.1 A More Flexible \hat{d}_M

In finite samples, p-values from tests of false null hypotheses will not always be zero. Similarly, many hypothesis tests are asymptotic and may not yield uniformly-distributed p-values in finite samples. Still, p-values from SSTs may exhibit something similar to the dichotomous behavior that motivates \hat{d}_M , in which p-values for $d \leq d^*$ are distributed differently than p-values for $d > d^*$. For this reason, ? suggested a more flexible estimate:

$$\hat{d}_M^{ab} \equiv \arg \min_{\hat{d} \in \mathbb{N}; 0 < b < a < 1} \sum_{d \leq \hat{d}} (p_d - a)^2 + \sum_{d > \hat{d}} (p_d - b)^2 \quad (3)$$

Like \hat{d}_M , model selector \hat{d}_M^{ab} looks for behavior that differs between p-values testing true and false null hypotheses. Unlike \hat{d}_M , it does not depend on theoretically established distributions for these p-values, but searches over a grid for their location parameters. \hat{d}_M^{ab} will be more computationally expensive to compute than \hat{d}_M , but will often yield better results, especially in small samples.

3.2 Edge Testing

Typically, the p-values from SSTs will be mutually correlated. This will be particularly pronounced in situations analogous to RDD bandwidth selection,

in which SSTs are used to choose among nested datasets. In this situation, a p-value for choice d , p_d , is based on the same data as the previous p-value, p_{d-1} , along with with sometimes only a few extra cases.

In contrast to specification tests that researchers use to check fully-specified models, and are designed to check the model as a whole, SSTs are an explicit and planned part of the model selection process. That being the case, their focus should be on differences between potential specifications, rather than on overall suitability. We refer to the former as “edge testing,” since it focuses hypothesis tests on edge cases, and the latter “total testing.”

When decisions d determine which data are included in the analysis, as in RDD bandwidth selection, the choice between edge and total testing is a choice between null hypotheses to test. The edge null is:

$$H_{0d}^{edge} : \mathcal{A} \text{ is true for } i \in \rangle_d \quad (4)$$

whereas the total null is

$$H_{0d}^{tot} : \mathcal{A} \text{ is true for } i \in \mathcal{I}_d \quad (5)$$

where, as above, $\mathcal{I}_d = \rangle_1 \cup \dots \rangle_d$, all data included in specification d .

For instance, in selecting a bandwidth for an RDD, as in ?, researchers test for, say, equality of means of a covariate x between treated subjects, with running variable values R at one side of the cutoff, and control subjects with R on the other side. Here d indexes candidate bandwidths, $\max |R - c| = bw_d$.

Then $\mathcal{I}_d = \{i : |R_i - c| = bw_d\}$ and $\mathcal{I} = \{i : |R_i - c| \leq bw_d\}$. Therefore, $H_{0d}^{tot} : \mathbb{E}[x|0 < R - c \leq bw_d] = \mathbb{E}[x|-bw_d \leq R - c < 0]$ and $H_{0d}^{edge} : \mathbb{E}[x|R - c = bw_d] = \mathbb{E}[x|R - c = -bw_d]$.¹ For the sake of demonstration, say $\text{var}(x) = \sigma^2$. For $d \leq d^*$, $\mathbb{E}[x||R - c| = bw_d] = 0$ but for $d > d^*$ $\mathbb{E}[x|R - c = bw_d] = \tau$ and $\mathbb{E}[x|R - c = -bw_d] = -\tau$. Further, say there are $n_d = n_0$ at each possible bandwidth bw_d . For $d = d^* + 1$, testing H_{0d}^{tot} means comparing the means of two samples of size $(d^* + 1)n_0$, each with standard deviation $\sqrt{\sigma^2 + \tau^2(1 - 1/d^*)}$ and with means $\pm\tau/d^*$. On the other hand, a test of H_{0d}^{edge} compares the means of two smaller samples, each of size n_0 , with standard deviation σ and means $\pm\tau$. As long as $d^* > 1$, the power of a t-test for H_{0d}^{edge} will be greater than the power for H_{0d}^{tot} ,² better allowing a the SST procedure to distinguish between d^* and d .

4 A Simulation Study

5 Two Data Examples

To illustrate these ideas—edge testing, and the change point and hypothesis testing approaches to selecting d —we will briefly illustrate them with two data examples. The two examples correspond to the two broad categories of specification we have discussed: selecting data to analyze and selecting a

¹These are simplifications of Assumption 4 in ?, which treats x as fixed, not random.

²The non-centrality parameter in the H_{0d}^{tot} test is $\frac{2\tau/d^*}{\sqrt{\sigma^2 + \tau^2(1 - 1/d^*)}} \frac{\sqrt{(d^* + 1)n_0}}{\sqrt{2}}$ and the non-centrality parameter in the H_{0d}^{edge} test is $\frac{2\tau}{\sigma} \frac{\sqrt{n_0}}{\sqrt{2}}$

model specification.

5.1 SSTs in Regression Discontinuity Bandwidth Selection: Estimating the Effect of Academic Probation on College GPAs

At many universities, students who fail to achieve a minimum GPA are put on academic probation (AP) (See, e.g. ?). This provides them access to a set of resources designed to address personal issues that may be hindering their performance. Perhaps more importantly, AP is a threat—students on AP who do not improve are subject to disciplinary measures such as suspension. ? recognized that AP can form a regression discontinuity design (RDD), in which treatment is a function of a “running variable” with a pre-determined cutoff. Specifically the treatment Z , students’ AP status, is (almost) a deterministic function of a “running variable” R , students’ grade-point-averages (GPAs). Students with a GPA below a pre-determined cutoff, $R < c$, are put on AP. That being the case, students with GPAs just below c may be comparable to students with GPAs just above c —comparing these two sets of students allows researchers to estimate the effect of AP on outcomes Y . The challenge becomes defining “just above” and “just below”; SSTs may be able to play a role here.

For example, ? (CFT) suggests directly comparing the outcomes of subjects with R very close to c , say with $R \in [c - bw, c + bw]$ for some bandwidth

$bw > 0$ To choose bw , CFT uses pre-treatment covariates \mathbf{X} , and covariate balance tests range of candidate bandwidths. For each possible bw , they test the hypothesis that the covariates are balanced:

$$\mathbf{X} \perp\!\!\!\perp Z | R \in [c - bw, c + bw] \quad (6)$$

and choose the largest bandwidth in which (??) cannot be rejected.

Bandwidth selection for RDDs, and the role of covariate balance tests, encompasses a growing literature. As its name suggests, regression discontinuity typically relies on regression modeling: the goal is to model Y as a function of R on either side of c to estimate the average treatment effect for subjects with R in an infinitesimally-small interval around the cutoff c (See ?). In contrast, CFT dispenses with regression altogether. One popular way to ensure robustness to model misspecification is to fit the regression models to a subset of the data with R in a window around c . A number of methods exist to choose an optimal bandwidth bw —the width of the window—that is both large enough to allow for precise effect estimation but small enough to ensure robustness. ? suggest using non-parametric estimates of the curvature of the regression function of Y on R , combined with local linear regression, to choose a bw that minimizes mean-squared-error. However, other authors have suggested choosing bw (or an analogous quantity) based on SSTs, including ?, which presents a Bayesian approach analogous to CFT’s, ?, which discusses the use of robust regression models, and ?, which proposes a method

to estimate effects for subjects with R farther from c . In the latter paper, SSTs do not test covariate balance, but the irrelevance of R conditional on covariates X , for subjects in a given bandwidth.

This section will illustrate several approaches to SSTs in the context of estimating the effect of AP for first year college students on subsequent GPAs. For the sake of simplicity, the discussion will be limited to CFT’s general approach to regression discontinuity designs; however, many of the SST methods can be extended to other RDD analyses. In their analysis, CFT considered a set of seven covariates: students’ high-school GPA (expressed in percentiles), age at college matriculation, number of attempted credits, gender, native language (English or other), birth place (North America or other) and university campus (the university consisted of three campuses). A version of Hotellings T^2 test that models treatment assignment Z , and not X , as random (?) is used to test balance. The resulting p-values, Total Testing—testing hypotheses (6)—and Edge Testing, testing $H_{0bw} : X \perp\!\!\!\perp Z || R| = bw$, are plotted in Figure ??.

In this case, total and edge testing paint similar pictures: covariates are imbalanced for most bandwidths. Both \hat{d}_M and $\hat{d}_M^{a,b}$, marked with a green vertical lines in Figure ??, select the lowest possible bandwidth. Since the first p-value 0.015 is below any α considered, a strategy that chooses the last non-rejected bandwidth (denoted above as $\hat{d}_{\hat{\alpha}}$) rejects every bandwidth. According to these methods, the CFT method is unsuitable for this dataset. However, the scattered large p-values at some bandwidths lead versions of \hat{d}_{α}

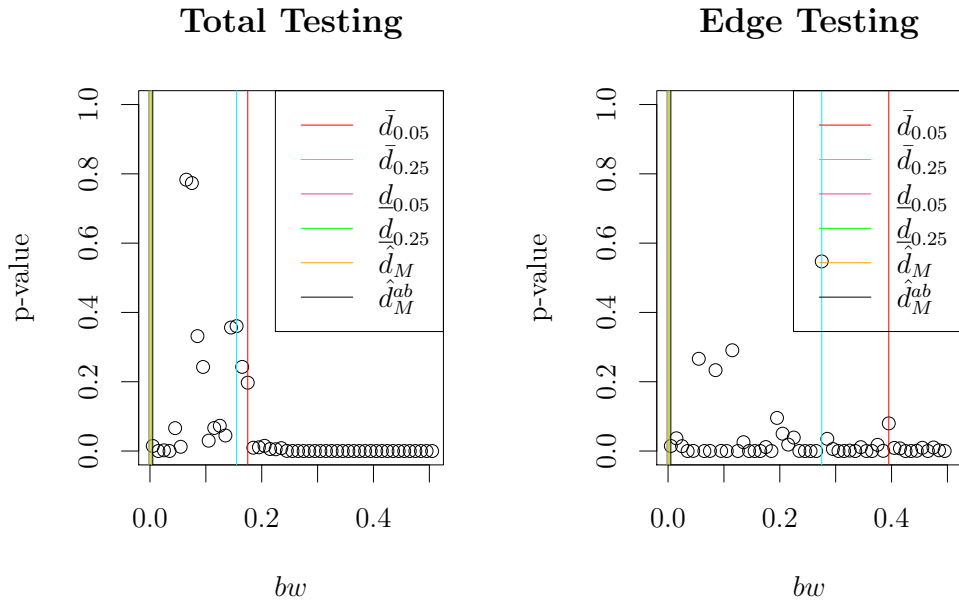


Figure 2: P-values from total and edge testing for balance in all seven covariates from the LSO analysis. Vertical lines denote bandwidth choices using different criteria.

Total Testing: High School GPA Edge Testing: High School GPA

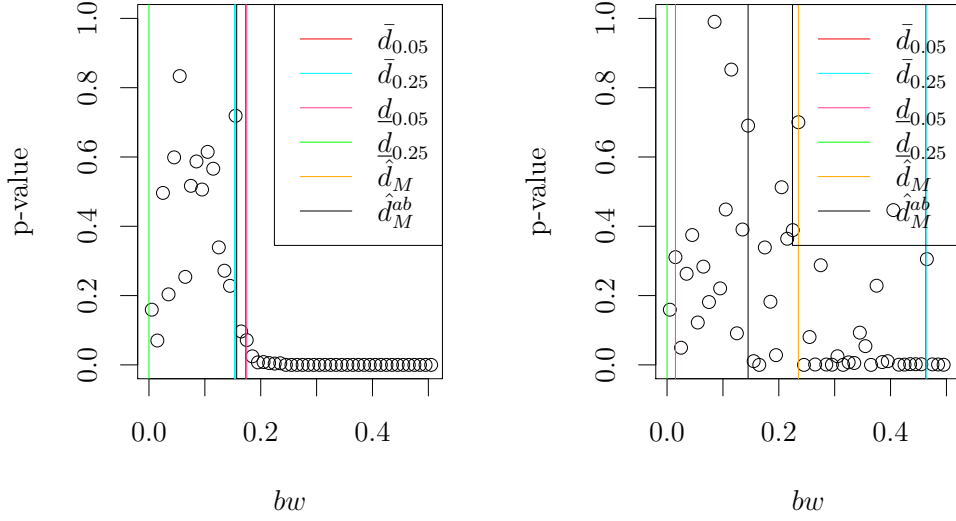


Figure 3: P-values from total and edge testing for balance in all seven covariates from the LSO analysis. Vertical lines denote bandwidth choices using different criteria.

to select larger bandwidths. This fact illustrates a weakness of \hat{d}_α : strong evidence against a model specification will be discarded in the presence of one favorable p-value.

To better illustrate differences between the window selection strategies, we consider the covariate high school GPA alone. Since the outcome of interest is itself a GPA, prior measures of GPA are arguably the most relevant and important to control. P-values from total and edge tests of balance in high school GPA are displayed in Figure 3. Fortunately for the illustration here, high school GPA may be balanced for small *bws*. P-values from total

and edge testing in 3 are markedly different. Edge testing shows some bandwidths, considered in isolation, appear to plausibly satisfy covariate balance, while others do not. On the other hand, the p-values from total testing are more nearly monotonic: imbalance in high school GPA at lower bandwidths causes balance tests to reject at higher bandwidths as well. The opposite also occurs: for instance, at bandwidth 0.155, total testing gives 0.719, whereas edge testing gives 0.011. The explanation is that high school GPA at that bandwidth is indeed imbalanced, but in the opposite direction of other imbalances: the AP students at $bw = 0.155$ had *higher* high school GPAs than those who were not on AP, to an extent that counteracted imbalances at smaller bandwidths.

That said, edge testing suggested larger bandwidths for almost all procedures.

This example suggests that edge testing might only be suitable for change-point-based window selectors. Further, the more flexible \hat{d}_M^{ab} outperformed the asymptotic estimator $\hat{d}_M^{0.5,1}$ by choosing $b = NA$, so that even moderately small p-values suggested departures from covariate balance.

5.2 Lag Order in AR(p) Models: US Total Unemployment

Figure ?? shows the United States total unemployment rate from 1890 to 2016. Assume that it follows an “AR(p)” model; that is,

$$unemp_t = \mu + \sum_{i=1}^p \phi_i unemp_{t-i} + \epsilon_t \quad (7)$$

where μ and $\{\phi_i\}_{i=1}^p$ are parameters to be estimated and ϵ_t is white noise. In this model, the unemployment in one year is a function of unemployment rates in the previous p years, but conditionally independent of even earlier measurements. More generally, we may write (7) as

$$unemp_t = \mu + \sum_{i=1}^{\infty} \phi_i unemp_{t-i} + \epsilon_t \quad (8)$$

with $\phi_i = 0$ for $i > p$.

Having settled on model (7), the analyst must choose p , the lag order. SSTs can be useful here (e.g. ?). Consider the null hypothesis $H_p : \phi_i = 0$ for all $i > p$; a reasearcher could test a sequence of such null hypotheses, for a set of plausible values of p , and choose the p based on the results. Of course, there are other options for choosing p , including substantive theory or optimizing information criteria, like AIC or BIC ((??) though ? points out that differences in AIC or BIC are essentially likelihood ratio test statistics). In the absence of substantive theory, SSTs can assist a modler to choose the

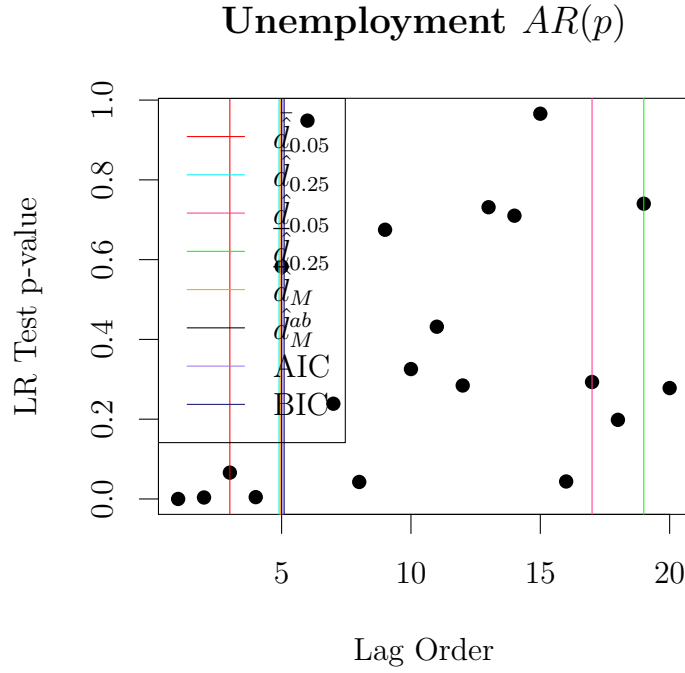


Figure 4: P-values from sequential likelihood-ratio tests of model fit, comparing models $AR(p)$ with $AR(p+1)$ in the annual total US unemployment rate (logged) time series.

smallest model that is still approximately correct—as opposed to the model that maximizes predictive accuracy as measured by, say, mean squared error.

6 Discussion