How to Measure Legislative District Compactness If You Only Know it When You See it*

Aaron Kaufman[†] Gary King[‡] Mayya Komisarchik[§]

December 29, 2020

Abstract

To deter gerrymandering, many state constitutions require legislative districts to be "compact." Yet, the law offers few precise definitions other than "you know it when you see it," which effectively implies a common understanding of the concept. In contrast, academics have shown that compactness has multiple dimensions and have generated many conflicting measures. We hypothesize that both are correct — that compactness is complex and multidimensional, but a common understanding exists across people. We develop a survey to elicit this understanding, with high reliability (in data where the standard paired comparisons approach fails). We create a statistical model that predicts, with high accuracy, solely from the geometric features of the district, compactness evaluations by judges and public officials responsible for redistricting, among others. We also offer compactness data from our validated measure for 17,896 state legislative and congressional districts, as well as software to compute this measure from any district. Word count: 9987

Replication Materials: Data, code, and other information needed to replicate all analyses in this article are available on the American Journal of Political Science Dataverse at Kaufman, King, and Komisarchik, 2020.

^{*}Winner of the 2018 *Robert H. Durr Award* from the Midwest Political Science Association. Our thanks to Steve Ansolabehere, Fred Boehmke, Ryan Enos, Dan Gilbert, Jim Griener, Bernie Grofman, Andrew Ho, Dan Ho, James Honaker, Justin Levitt, Luke Miratrix, Max Palmer, Stephen Pettigrew, Jamie Saxon, Steve Shavell, Anton Strezhnev, Wendy Tam, Rocio Titiunik, Larry Tribe, Robert Ward, participants in "A Causal Lab", and the audiences at the Society for Political Methodology Meetings, Nuffield College at Oxford University, the Harvard Applied Statistics Workshop, and the ICPSR Summer Program for helpful data or suggestions; and to Stacy Bogan, the Center for Geographic Analysis, and the Institute for Quantitative Social Science at Harvard University for research assistance and support.

[†]Assistant Professor, New York University Abu Dhabi, Saadiyat Island, Abu Dhabi, UAE; AaronrKaufman.com; aaronkaufman@nyu.edu, (818) 263-5583.

[‡]Albert J. Weatherhead III University Professor, Institute for Quantitative Social Science, 1737 Cambridge Street, Harvard University, Cambridge MA 02138; GaryKing.org, King@Harvard.edu, (617) 500-7570.

[§]Assistant Professor, University of Rochester, Harkness Hall, 333 Hutchinson Road, Rochester, NY 14627; scholar.harvard.edu/mkomisarchik; mkomisarchik@fas.harvard.edu, (720) 220-9328.

1 Introduction

Compactness is treated in the law as an important legal bulwark against gerrymandering. The Apportionment Act of 1901, many court decisions, and 18 state constitutions require compactness for U.S. House districts, and 37 states require their legislative districts to be compact (see j.mp/aRED). Compactness is also required in federal law as one of the "traditional redistricting principles" which, when followed, can "defeat a claim that a district has been gerrymandered..." on the basis of race (*Shaw v. Reno*, 509 U.S. 630, 647, (1993)) or political party (*Davis v. Bandemer*, 478 U.S. 173, 2815, (1986)).¹

Compactness is also important for the academic literature, where scholars seek to help the redistricting and litigation processes, and also to study venerable political science questions such as the causes, consequences, and normative implications of compact districts over American history (e.g., Ansolabehere and Palmer, 2016; Ansolabehere and Snyder Jr, 2012; Forgette and Platt, 2005). Compactness intuitively refers to both how close a legislative district's boundaries are to its geographic center and how "regular" in shape a district appears to be. But upon deeper study, scholars have shown that in fact compactness is a complicated multidimensional concept and have offered almost 100 measures of different features of it (e.g., Niemi, Grofman, Carlucci, and Hofeller, 1990).²

While many state constitutions explicitly require compactness, the vast majority provide no definition or measure for how to detect violations of the standard. For example, the Constitution of Illinois says only "Legislative Districts shall be compact". The Constitution of Hawaii requires that "Insofar as practicable, districts shall be compact." In Arizona, the Constitution orders that "Districts shall be geographically compact and con-

¹Claims about most other types of unfairness in redistricting all also seem to depend on a legal finding of noncompactness (*Davis v. Bandemer*, 478 U.S. 165; Justice Powell in *Vieth v. Jubilerer*, 541 U.S. 267 (2004) 176-177; *Kirkpatrick v. Preisler*, supra, at 394 U. S. 526, 538).

²The empirical claim sometimes implied in the law, that compactness requirements constrain racial or partisan gerrymandering, is the subject of active research program (Altman and McDonald, 2012; Barabas and Jerit, 2004; Chen and Rodden, 2013), and the role of compactness in ensuring other important normative virtues — such as better knowledge, communication, and trust between representatives and citizens — is also contested (Cain, 1984; Pildes and Niemi, 1993). But regardless of the outcome of these important debates, the degree of compactness of legislative districts will always have an essential role in defining the nature of representation and electoral competition in modern democracies, and an accurate measurement is essential to each debate.

tiguous to the extent practicable."³

The federal courts have been similarly vague. They have acknowledged both the multitude of possible measures for compactness, and the fact that they often produce different conclusions.⁴ Except in rare cases, the courts have not provided guidance on particular measures or seen the need for them. For example, Justice Souter stated that "it is not necessary now to say exactly how a district court would balance a good showing on one of these indices against a poor showing on another, for that sort of detail is best worked out case by case" (*Vieth v. Jubelirer*, 541 U.S. 267 (2004); Souter dissenting). And most famously, a Supreme Court opinion indicated "One need not use Justice Stewart's classic definition of obscenity—'I know it when I see it'—as an ultimate standard for judging the constitutionality of a gerrymander to recognize that dramatically irregular shapes may have sufficient probative force to call for an explanation" (*Karcher v. Daggett*, 462 U.S. 725, 755 (1983)). Here, the Court at once laments the absence of a single quantitative standard while also implying that the concept is clear enough that all reasonable observers should understand it in the same objective way.

Consistently invoking the idea of "compactness" without a clear definition or required measure suggests two conclusions about the law. First, the law seems to imply that "compactness" is a single, coherent, and agreed upon concept, discernable simply by examining a district map. After all, how could the courts expect legislators to draw districts that comply with "compactness" without a shared understanding of what it means? And second, this lack of precision in the law has given redistricters and litigants battling over legislative maps in specific cases wide latitude to choose measures of compactness and reach

³Some states have passed laws highlighting certain features of compactness that may help with intuition but neither precision nor application. For example, Virginia Senate Joint Resolution 224 (1/14/2015, Article II, Section 6(5)) reads "Each legislative and congressional district shall be composed of compact territory. Districts shall not be oddly shaped or have irregular or contorted boundaries, unless justified because the district adheres to political subdivision lines. Fingers or tendrils extending from a district core shall be avoided, as shall thin and elongated districts and districts with multiple core populations connected by thin strips of land or water...." Iowa (Iowa Code, Title II §42.4) and Michigan (Congressional Redistricting Act 221 of 1999, Redistricting plan guidelines) mention some precise measures but not how to use this information.

⁴"Indeed," writes Justice Souter, dissenting in *Vieth v. Jubelirer*, "although compactness is at first blush the least likely of these [traditional redistricting] principles to yield precision, it can be measured quantitatively in terms of dispersion, perimeter, and population ratios, and the development of standards would thus be possible."

opposing conclusions (Defendant-Intervenors' Post-Trial Brief [at pp. 18], *Bethune-Hill v. Va. State Bd. of Elections*, 141 F. Supp. 3d 505 (E.D. Va. 2015) (No. 3:14 Civ. 852), ECF No. 104; and Motion In Limine Regarding Plaintiffs' New Compactness Test [at pp. 4], *Vesilind v. Va. State Bd. of Elections*, No. CL 15-3886 (Va. Cir. Ct. 3/31/2017).). Even when litigants might agree on the compactness of any one district by knowing it when they see it, systematically judging the compactness of many districts, or an entire redistricting plan, cannot be accomplished by merely looking. As such, the courts and policy makers tend to get very little benefit from the decades of work on quantitative measures of compactness offered by social scientists.

We attempt to span this divide between the seemingly universal understanding of compactness proposed in or needed for the application of the law, and the theoretical complexity and multidimensionality revealed in the social science literature. We do this by inferring, measuring, and validating the single underlying dimension of compactness that practitioners may need to apply the law, and we find that people of all types seem to agree upon it. In other words, since compactness in the law is, for all practical purposes, defined by the judgment of human observers — including redistricters, experts, consultants, lawyers, judges, public officials, and ordinary citizens — the claim of an objective standard, measured on a single dimension, can only be supported if most educated people evaluated a district's compactness in the same way. We provide this objective measure and show that these and other groups of observers all view compactness in accordance with it. This new dimension is not the average (or principal component) of existing measures but a new quantitative construction that accurately and reliably predicts human judgment.

In four sections, we proceed by *conceptualizing*, *measuring*, *validating*, and *interpret-ing* our derived dimension of compactness. Section 2 inductively defines the underlying dimension by building on the encyclopedia of existing diverse measures, adding new ones that show how humans perceive objects like district shapes, and providing intuition about the commonly perceived dimension we seek to measure. Section 3 then develops a way to measure this concept by eliciting views of the compactness of specific districts from respondents using a novel survey approach to rank order districts according to their com-

pactness. We are forced to develop a new method because the standard approach in the survey literature to a problem like this, Thurstone's paired comparisons, completely fails in our application. The high levels of intercoder and intracoder reliability produced by our alternative approach are consistent with a unidimensionality hypothesis (and suggests that our survey methodology may have other applications). This section then uses these results to build a statistical model that predicts with high accuracy how individuals rank districts, given only the the districts' shapes.

Our results enable us to apply one of the most important principles of statistics — defining the quantity of interest separately from the measure used to estimate it — and, as a result, to provide evaluations that make our approach vulnerable to being proven wrong. We do this in Section 4 with cross-validation and then extensive out-of-sample validations in samples of public officials and judges from many jurisdictions, as well as redistricting consultants and expert witnesses, law professors, law students, graduate students, undergraduates, ordinary citizens, and Mechanical Turk workers. Application of this same principle also enables us to provide the first uncertainty estimates for a measure of compactness offered in the literature (see Supplementary Appendix D). Section 5 then offers interpretations of the resulting measure, and Section 6 concludes.

2 Conceptualizing

We now attempt to inductively characterize the concept of compactness that most laws, constitutions, judicial opinions, and participants in redistricting at least implicitly assume human observers intuitively understand.

As districting is "one area in which appearances do matter" (*Shaw v. Reno*, 509 U.S. 630, 647, 1993), our approach is to measure the absolute compactness of the geometric shape of a district, separately from other facts that can impact this measurement such as geography or population. This is the most common basis for a compactness definition, dating well before the famous "Gerry-Mander" cartoon (Tisdale, 1812), but not the only one possible. Absolute compactness, in turn, may be constrained or influenced by fixed features of the state geography, such as rivers, coastlines, or highways. We measure the

quantity that would be influenced by these features, so that it measures the concept in the law and can be useful for further research. If a researcher had the alternative goal of defining and measuring relative compactness, based on how close it is to a realistic ideal, then our measure would be a key component in that calculation. In addition to measuring absolute compactness based on shape, our methods can also be used to measure compactness based on other criteria, such as population dispersion (Fryer Jr and Holden, 2011; Hofeller and Grofman, 1990; Niemi, Grofman, Carlucci, and Hofeller, 1990); see Section 3.3.

We attempt to characterize the compactness of each district separately. Although changing the boundaries of one district obviously affects neighboring districts, separate measurement follows major redistricting litigation, which typically evaluates the compactness of districts individually or in a small group rather than for an entire state redistricting plan all at once (e.g., *Shaw v. Reno*, 509 U.S. 630 (1993), pp. 637, 647, 656). This strategy is especially useful for the most fine grained scholarly research on the causes and consequences of compactness.⁵

Section 2.1 highlights empirical inconsistencies in existing shape-based measures to convey that the possible conceptual definitions of compactness, underlying these measures, are multidimensional. Then Section 2.2 provides intuition and tools to build toward a single concept of compactness.

2.1 Multiple Dimensions Underlying Existing Measures

Numerous specific compactness measures have been proposed in the academic literature, each one fitting different qualitative conceptual definitions and intuitions for certain ge-

⁵Aspects of the overall methodology we develop here can also be applied to some other redistricting criteria, when additional data are available (or to concepts unrelated to redistricting that you only know when you see). These may include other characteristics of districts such as size; population equality across districts; where people live within a district (Fryer Jr and Holden, 2011); whether the district divides communities of interest or local political subdivisions; whether incumbents are paired or grouped in the same district and so have to run against each other to keep their jobs; what types of people are included in or excluded from a district; and, as a result, partisan fairness, electoral responsiveness (Gelman and King, 1994; Grofman and King, 2007), and racial fairness (King, Bruce, and Gelman, 1996). Redistricting also influences more personalistic factors common in real redistricting cases, such as whether a specific district includes features like a military base (which can influence a candidate's policy preferences) or a prison (which counts under "equal population" requirements but not votes), or even whether a candidate's parents homes or children's schools are drawn out of his or her district.

ographical configurations and violating it for others (Altman, 1998; Niemi, Grofman, Carlucci, and Hofeller, 1990; Stoddart, 1965; Young, 1988). These measures are based on geometric concepts such as perimeters, areas, vertices, and centroids, often in comparison with some pure form geometric object such as a circle, rectangle, polygon, or convex hull. Each, however, focuses on a different dimension of what might be called compactness. Consider, for example, the five most frequently used measures by academic researchers, and also by experts in redistricting litigation: Length-Width Ratio, the ratio of the length to the width of the minimum bounding rectangle (Harris 1964; Timmerman, 100 N.Y.S. 57, 51 Misc. Rep. 192 (N.Y. Sup. 1906)); Convex Hull, the ratio of the area of the district to the area of the minimum bounding convex hull; Reock, the ratio of the area of the district to the area of a minimum bounding circle (Reock, 1961); Polsby-Popper, the ratio of the area of the district to the area of the circle with the same perimeter as the district (Polsby and Popper, 1991; Schwartzberg, 1965); and (modified) Boyce-Clark, the (normalized) mean absolute deviation in the radial lines from the centroid of the district to its vertices (Boyce and Clark, 1964; Kaiser, 1966; MacEachren, 1985). For details on these and others, see Supplementary Appendix A.

Without a gold standard, we cannot determine any measure's formal statistical properties, its error rates, or when it might fail. Although different measures are sometimes correlated, choices among these are presently made by qualitative judgment. Creative scholars have managed to use existing measures productively in research by combining multiple measures, adjusting or weighting each for specific purposes, or making careful qualitative decisions in specific cases (Ansolabehere and Palmer, 2016; Niemi, Grofman, Carlucci, and Hofeller, 1990).

We illustrate the issues with measuring compactness by presenting Figure 1, four state house districts from Alabama in 2000. Readers may wish to draw their own conclusions about the relative compactness of these districts, but we now provide in Table 1 an indication of how the most popular five measures rank them (we discuss X-Symmetry and significant corners in Section 2.2). As can be seen from the first five rows of Table 1, every one of these measures gives a different rank order for the four districts. We introduce

two new compactness measures in Section 2.2 for a different purpose; these are given at the bottom of Table 1 and also give unique rankings of the same districts. This example is merely a proof of concept, but finding such examples is easy: By random sampling, we estimate that in our collection of 17,896 state legislative and congressional districts (see Supplementary Appendix B), there exist 162 trillion sets of four districts such that every one of the seven measures provides a unique rank order. Of course, there is a large number from which to choose (this large number being about 0.15% of the total), but inconsistencies among in rankings on fewer than seven measures is both commonplace and is congruent with the long literature on this subject.

[Figure 1 about here.]

[Table 1 about here.]

2.2 Toward a Single Compactness Dimension

We now provide intuition helpful in turning the multiple types and dimensions of compactness illustrated in Section 2.1 into a single unidimensional concept underlying common conceptions, but in the absence of political or personal biases. We continue to proceed inductively, with Section 3 devoted to measuring this concept. We do this in three ways, followed by a characterization of the dimension of interest.

First, our goal is to elicit views about compactness, but without the biases psychologists have long demonstrated skew human judgments in the direction of our own political and other preferences. Although such unbiased views may be the goal of lawyers advocating on behalf of their clients, research has shown that subject matter experts are as vulnerable to bias as nonexperts, and more overconfident in the belief that they can avoid it. The only reliable solution has been to remove even the possibility of bias by instituting formal procedures (such as double blind experiments). (See Kahneman, 2011). We thus elicit views about compactness without revealing to respondents how their decisions in any one situation might benefit one political party or another. This is a critical point: Because individual judges, advocates, redistricters, and experts do not have access to the mental processes in their own thinking that would enable them to evaluate and avoid these biases (Wilson and Brekke, 1994), they would also be unable to come to the same judgment as our measure in the context of a real redistricting contest by merely looking at a district shape.

Second, all existing compactness measures are *rotationally invariant*, meaning that if we rotate a district, say 45 degrees, a measure will have the same value. Although this is a reasonable normative standard from some perspectives — and we discuss below how to easily adjust our methods to impose this restriction if desired — human beings (including judges) do not evaluate districts in this way. In fact, human perception is famously sensitive to the rotation of objects: even familiar faces can become unrecognizable when viewed upside down (e.g., Maurer, Le Grand, and Mondloch, 2002). Our own experimentation done in R Shiny (Kaufman, 2020) suggests that people view long thin district shapes located on a diagonal () as less compact than the same shape located along the horizontal axis ().⁶ In contrast, legislative districts always have a well defined up (north) and down (south), as displayed on every commonly used map. Indeed, courts, redistricters, and judges virtually always use this single standard orientation and do not rotate districts when evaluating compactness; as a result, their decisions are not rotationally invariant. In other words, since the usual orientation of a district has precedence in how humans interpret it, some of our measures need to pick up on these features.⁷

Thus, primarily for illustration in this section, and later as a measurable feature of district shape that can be included (and if desired controlled) in our statistical model, we define here a new compactness measure that is not rotationally invariant. We do not intend this measure to substitute for other measures or to even be especially important on its own, but it will be useful to represent human perception. Thus, we define *X-Symmetry* by dividing the overlapping area, between a district and its reflection across the horizontal axis, by the area of the original district. Shapes like circles and rectangles have overlap regions equal to that of the original district and so have X-Symmetry values of 1. A long

⁶This pattern may be related to the "horizontal-vertical illusion" discovered in psychology (Prinzmetal and Gettleman, 1993).

⁷We note as well that, since all modern political boundaries are drawn with respect to cardinal directions, those directions are necessarily considered in examining districts.

thin district stretched out from top left to bottom right, or one like , have X-Symmetry values close to zero. This measure, applied to the four districts in Figure 1, gives unique rankings for each; see the sixth row of Table 1.

Since we are attempting to quantify human perception, we try to avoid imposing theoretical notions of what compactness should be, what might be rational, or what meets various mathematically "pure" standards that implicate one normative preference or another (such as rotational invariance). Finding the common objective measure that exists in minds of districting authorities, the courts, and others requires respecting how humans think rather replacing it with alternative normative preferences. Although the courts have never addressed the issue, in all likelihood those who drafted compactness requirements in legislative statues, judicial opinions, and state constitutions, that imply that the concept is so simple that you know it when you see it, were not assuming rotational invariance. However, if a rotational invariant measure is desirable or at some point required, we can easily impose it using a procedure analogous to what we do for avoiding political bias. Thus, we would use all the procedures described in this paper except that we would simply display districts at random rotational angles when eliciting compactness evaluations.

Third, another feature of human perception is how we define what constitutes a "significant" feature of a district. If a roughly circular district has a ragged border, which of the small border inlets and peninsulas count as notable deviations from the circular shape? For example, suppose we give a large number of people the task of drawing from memory the shape of the continental United States. These drawings will all differ, but they will likely all include some of the same features — a roughly rectangular shape, a peninsula for Florida, a larger one for New England, and perhaps a somewhat rounded western ocean boarder. In other words, despite the enormous number of specific small features and vertices along the boarder to choose from, virtually all Americans are likely to recall, thus judging as significant, a small number of the same features.

To include this highly qualitative feature of human perception, we consider algorithms computer scientists design to list all of the "objects" in an image. There is no correct answer, but it turns out that different people are likely to give similar answers, and the automation goal is to list the objects a human would identify. As we do with X-Symmetry, we illustrate this idea quantitatively, and give an example that will later become part of our model. To do this, we turn the geometric district shape into a set of pixels (i.e., changing from vector to raster representation), apply a corner detection algorithm (Shi and Tomasi, 1993), and count the number of "significant" corners. The more significant corners, the less compact the district by this metric. The last row of Table 1 gives the rankings of the four districts in Figure 1 according to the number of significant corners. This measure also gives the four districts a unique ordering.

Finally, we try to convey intuition about the underlying dimension of compactness we will quantify in the next section. We do this visually, by presenting in Figure 2 a set of districts that range from most (panel a) to least (panel d) compact. We find that almost anyone familiar with the district-based nature of modern democracy, and some sense of the word compactness, finds that district (a) is more compact than (b), which is more compact than (c), which is more compact than (d). The question is how to quantify this notion, so that it works for these four districts and all other geometric shapes, a topic to which we now turn.

[Figure 2 about here.]

3 Measuring

We now develop a more explicit measure of the concept of compactness to satisfy our requirements in Section 2. The immediate quantitative goal of the procedure is a *continuous* measure for each district, between 1 and 100, that estimates the expected rank a respondent would assign a district if embedded in a set with 99 others. With this measure, we can rank order any set of n districts, given only quantitative measures of their geometric shapes.

To construct this measure, we first develop a method of eliciting views about compactness directly from survey respondents, something universally recognized as important but rarely done in this literature except informally by researchers (Angel and Parent, 2011; Chou, Kimbrough, Murphy, Sullivan-Fedock, and Woodard, 2014). Appendix A attempts this by applying best current practices in survey research — using a modern version (David, 1988) of Thurstone's venerable paired comparisons (Thurstone, 1927), a method that dates at least to 1860 (Fechner, 1966). Under this approach, we pose a set of simple survey questions, each asking the respondent to decide which of two districts is more compact and, from the many answers, we construct the full ranking. We explain the motivation behind this approach and then demonstrate empirically that it *utterly* fails to accomplish its goal for this application.

Given the failure of paired comparisons, we have no choice but to develop a new approach. Thus, in Section 3.1, we turn to the method that paired comparisons was originally designed to supplant — asking respondents to rank many districts all at once. We show that, as we apply it, this approach turns out to work extremely well in our application (and may also work for many others too). As we describe, the supposed advantages of paired comparisons turn out to be disadvantages and the disadvantages of ranking turn out to be advantages. Section 3.2 takes the resulting survey elicitation method as our outcome variable, and new gold standard, and builds a statistical model to predict it from geometric features of the districts. Details about data used appear in Supplementary Appendix B.

3.1 How Ranking Outranks Paired Comparisons

Why does the method of paired comparisons perform so poorly? We propose four reasons, which together leads us to a workable approach for our application, full ranking — the method which paired comparisons originally supplanted.

First, although $\binom{n}{2}$ paired comparisons is vastly smaller than n! rankings (see the start of Appendix A), for some applications rankings make be quicker. After all, how long would it take to carefully and accurately rank 20 district shapes by their degree of compactness (or 20 friends by their heights or 20 animals by their friendliness)? A lot less than 2 quintillion seconds. What the idea behind paired comparisons seems to miss is that humans are excellent at pattern recognition and seeing the big picture. Humans also intuitively apply time-saving heuristics that reduce the complexity of tasks, such as in our application by grouping districts into distinct types, and considering all members of the group at once before analyzing members within the group.

Thus, in practice with full ranking, we have tried to ensure that respondents are using their big picture skills, such as by suggesting to them that they simplify the task by working hierarchically, first grouping districts into three coarse groups, and then producing groupings within each group, and finally starting from the top and checking and adjusting each district's position within the ranking; however, we found that heuristics and intuitions are strong enough that dropping these instructions did not degrade our full ranking approach. We also tried full ranking with districts printed on paper and arrayed on a long table, as well as via an online system we built that allows districts to be dragged and dropped to their chosen location; we find no evidence that the mode of administration matters either (as with Blasius, 2012).

Second, human respondents work better when motivated and engaged. While paired comparisons successfully avoid the risk of asking respondents questions they do not understand, it is also an unavoidably boring and tedious task, especially after the first few questions. In contrast, ranking a large set of districts is more intellectually challenging and engaging (Fabbris, 2013). Our own cognitive debriefing strongly supports the advantages of ranking in this regard.⁸

Third, if it is possible for a survey respondent to rank (say) 20 districts without much trouble, then we can save considerable time by administering this one engaging survey task rather than having to ask 190 tedious paired comparisons for each respondent. Ranking would then save considerable time, expense, and respondent fatigue (Ip, Kwan, and Chiu, 2007). As a hint that this might work, Krosnick (1999) (studying rating rather than paired comparisons) finds that often "rankings give higher quality data than ratings".

And finally, the literature makes clear that compactness is a multidimensional concept (Niemi, Grofman, Carlucci, and Hofeller, 1990). Yet, we are trying to tap into a single unidimensional concept of compactness that we hypothesize respondents, if given the choice, would select and use. In this light, the fact that Thurstone's approach enables respondents to make each paired comparison *independently* of the others allows, and may

⁸We also experimented with having two coders participate together in ranking each set of districts, on the theory that the social connections would make the task even more engaging. Our results support this theory, in that respondents spent about 30% more time together completing the task, but this engagement was unnecessary since it did not increase inter- or intracoder reliability.

even encourage, them to use different dimensions for different comparisons. In other words, while "roundness" may be the deciding factor for compactness in one given pair of districts, length vs. width may be the relevant question in the next pair, and so forth. This may then be what results in the low levels of intercoder and intracoder reliability we have documented. In contrast, ranking has the advantage of encouraging respondents to *choose* a single dimension of compactness and to use it for all their decisions. With paired comparisons, the only way to do this would be to ask respondents to choose a single dimension explicitly and to keep that dimension in their heads while they answer 190 randomly ordered survey questions. Although the goal of any survey question is to be clear enough so respondents are answering the question intended by the researcher (i.e., on the dimension of interest), giving respondents multiple separate questions makes this difficult to achieve.

To test our hypothesis that ranking will work better than paired comparisons, we ask respondents to give a full rank order for 100 separate legislative districts by their degree of compactness.

To begin, we embed our 40 districts (which we used in 20 pairs in the experiments in Figures 7 and 8) among 60 others and ask a new set of respondents to rank all 100. To compute a relative assessment of the two methods, we evaluated intercoder and intracoder reliability of the *implied* paired comparisons of how these 20 pairs were ordered by full ranking and compared them to reliability from the *actual* paired comparisons. That is, from full ranking, we record only which district in each pair of 20 comparisons is ranked higher. Then, to compute intracoder reliability, we waited two weeks, shuffled the rank ordering, and asked the same respondents to rank the same 100 districts, again only using the 20 designated pairs among these. We then computed the percent agreement over time in these implied paired comparisons exactly as we did for the actual paired comparisons. The results, which appear in the same two figures (salmon colored histogram, at the right of each figure), are far more clearly separated from the random placebo test and have much higher levels of intracoder reliability than the actual paired comparisons. For intercoder reliability, in Figure 7, we have 75% agreement on average, and for intracoder reliability.

in Figure 8, we have 88% agreement on average.

Now that we have a method that bests paired comparisons for measuring compactness with respect to pairwise intracoder and intercoder reliability, we turn to evaluating full ranking on its own terms. We begin with intercoder reliability by correlating the ranks for 100 districts coded independently by (all possible) pairs of respondents. We then present in Figure 3 one scatterplot representing the pair of coders with the median correlation ($\rho = 0.77$ in the top left panel) as well as the pair with a correlation in the first quartile (bottom left) and in the third quartile (top right). In the bottom right of the same figure (salmon colored), we also present a density estimate (using a kernel truncated at the minimum and maximum observed correlations) of all the correlations, along with a baseline density estimate of correlations among randomly generated ranks. The conclusion from this figure reveals high intercoder reliability, clearly distinguishable from chance, and with no systematic error patterns in any individual scatterplot.

[Figure 3 about here.]

We then repeat this process for intracoder reliability by correlating the ranks for each respondent with the same respondent, re-ranking the same districts, two weeks later. Figure 4 shows these results in the same format as Figure 3. As would be expected, our results here are even stronger than for intercoder reliability. The median correlation (top left) is $\rho = 0.9$, with not much spread around the median (see salmon colored histogram in the bottom right panel). None of the scatterplots show any systematic patterns in deviations from the 45° line, and all indicate high levels of intracoder reliability.

[Figure 4 about here.]

3.2 A Statistical Measurement Model

To construct our ultimate measure of compactness, we begin with a set of districts and elicit the views of respondents via our full ranking survey approach. In Section 3.1, we describe this survey methodology. Supplementary Appendix B gives details of how we recruited our survey respondents, collected our set of districts, conducted our experiments,

and wrote and presented the ranking task to respondents. We also discuss there the mechanism for how we elicited ranking preferences, both in person (sorting paper cards with districts printed) and online (dragging and dropping district images).

Our data collection process results in six sets of 100 districts, each ranked by several individuals or pairs of individuals working independently. We average away random error by calculating the first principal component of the rankings produced for each set of 100 districts, preserving the ranked scale. This first principle component, a summary of human-derived compactness rankings, forms the outcome variable in our statistical model, using only information from the shape of districts as predictors. To produce our predictor variables, we calculate a set of geometric features including all seven compactness indicators from Table 1 and the others described mathematically in Supplementary Appendix A.

Finally, we train an ensemble of predictive methods with these data, consisting of least squares, AdaBoosted decision trees, support vector machines, and random forests. Supplementary Appendix C gives the details of these methods and of how we construct this ensemble and its component parts.

All further details and code are available in our replication data file which accompanies this paper. In the same way that logit or ordered probit take discrete outcome variables and generates continuous predictors, our training data consists of integers from 1 to 100, but our ensemble model produces continuous outputs.

3.3 Compactness as Shape and Population Dispersion

As described in Section 2, the concept of compactness in the law, most of the literature, and our paper is based on district shape alone. However, other conceptualizations may be of interest for some purposes, such as based on population, communities of interest, not dividing political subdivisions, etc. For each of these, all the methodological procedures we developed in this paper should be directly applicable. The measure that results from the application of our procedures entirely depends, of course, on the quantity of interest being estimated, and there is no guarantee that a measure of compactness based on shape will be related to one based on other criteria.

As one small proof-of-concept of the applicability of our approach, we repeated our survey with district shapes that also represented where in each district people live (An-solabehere and Palmer, 2016; Niemi, Grofman, Carlucci, and Hofeller, 1990). We ran this population distribution experiment with six undergraduates from different universities on the same set of 20 districts. Details of the experimental protocol appear in our replication data set. Results indicate that the median correlation between the $\binom{6}{2} = 15$ possible pairs of rankings was a substantial 0.94, with a range of 0.12. This is comparable to the results we found using shape alone.

4 Validating

Via cross-validation (in Section 4.1) and out-of-sample prediction in diverse populations (in Section 4.2), we now evaluate our single, unidimensional compactness measure, deterministically computed from a district shape, and confirm our hypothesis that the theoretical concept we are measuring is the same one people know when they see. The data for this section come from diverse populations including participants directly involved in decision making about legislative redistricting.

4.1 Cross-validation

We evaluate our model here with cross-validation, where each fold reserves one of our six sets of 100 districts. To do this, we use six groups of survey respondents, potentially making it harder for our model by mixing size of group, mode of administration, and type of respondent: (1) two pairs of undergraduates (the two within each pair working together) and one pair of graduate students; (2) one pair of undergraduates, one individual undergraduate, and one pair of graduate students; (3) 5 individual undergraduates, 5 pairs of undergraduates, and 16 Mechanical Turk workers; (4) 5 pairs and five individual undergraduates; (5) 8 undergraduates; (6) 8 undergraduates. (We found ex post that respondents gave similar rankings regardless of whether they worked alone or in pairs. Similarly, Mechanical Turk workers, undergraduates, and graduate students gave similar rankings on the same sets of districts.)

We then trained our model on groups 1–5 of respondents taken together, and predicted the remaining "test set" of respondents in group 6; we repeated this six times in total, with each group taking its turn as the test set and the remaining groups as the training set. The prediction from this model uses all information from the training set but only the district geometry (i.e., no survey information) from the test set. Figure 5 evaluates the performance of this procedure by providing six scatterplots corresponding to each of our training set-based predictions (horizontally) by the true test set values (vertically). As is evident, these cross-validation results indicate very high predictive accuracy. Correlations between predictions and test set values range from 0.92 to 0.96, with no noticeable systematic error patterns in any graph.

[Figure 5 about here.]

4.2 Predictive Validation in Diverse Populations

The statistical model in Section 3.2 is designed to predict human judgment about the compactness of any set of districts, given only the geometric shapes of the districts. Our model can make a prediction for any legislative district shape, including new districts and those that do not appear in our training set.

Our hypothesis is that any informed human being will judge the compactness of a set of districts in almost the same way, thus admitting to high levels of statistical reliability. We now test this hypothesis by asking a wide range of groups to evaluate the compactness of different sets of legislative districts and comparing these evaluations to our predictions. Our main test comes from 96 sitting justices, judges, and public officials, all with some responsibility for redistricting or deciding redistricting cases. We also elicited the views of 102 others, ranging from less to more involved in and knowledgeable about redistricting, including Mechanical Turk workers, who received small monetary payments, undergraduates, some of whom received hourly wages, and others who were not paid, including political science PhD students, law students, law faculty, redistricting consultants and expert witnesses, and lawyers involved in legislative redistricting cases.

We promised our respondents confidentiality, including their responses and the fact of

their participation. This was most obviously a concern in recruiting judges and justices, who decide redistricting cases, and other public officials, who have decision making authority in or substantial influence on the process. It turned out to be of no less a concern for some lawyers who try redistricting cases, and some consultants and expert witnesses who are held to account for their previous statements and opinions. For these reasons, we are not able to make these data available publicly, although we do make available the software we designed to let respondents sort districts online and all our specific experimental protocols. All these steps were approved by our university Institutional Review Board. (We have also prepared and field tested teaching exercises for American government classes that use our districts, enable students do the ranking exercise themselves, and compare them to our predictions.)

In this experiment, we asked each respondent to rank order twenty legislative districts, not included in our training data, by their degree of compactness; we represent the degree of predictive accuracy by a simple correlation with our predictions. All respondents ranked the same twenty districts. We portray our results in Figure 6 with a histogram for each of nine categories of people. As a baseline, we present a density estimate (in blue) of the percent agreement among random rankings, which is of course centered at zero, and the variance of which conveys uncertainty given n = 20 districts. The (salmon-colored) histogram is for Mechanical Turk workers. The remaining histograms of correlations appear in white, with black outlines. We do not distinguish among these for a further level of confidentiality, but they all lead to the same conclusion of very high levels of predictive accuracy.

[Figure 6 about here.]

We found no statistically significant differences between the size of the correlations among different groups of respondents. The main predictor of the strength of the correlations was the time spent on the task, with longer times yielding higher correlations. This accounts for the larger variance of Mechanical Turk workers, as they are paid by the completed task regardless of how long they spend.

18

5 Interpreting

Having conceptualized, measured, and validated our estimate of compactness, we now interpret the result. Of course, we already have one interpretation — that we know it when we see it. That is, our fully automated quantification of the compactness of a district geography reproduces how informed human observers evaluate a never-before-seen district shape. Our model can do this quickly for millions of potential districts in ways no human could ever do — and so it could be used in a court case comparing entire legislative plans or in academic research comparing many legislatures — but the quantity being estimated by our model and by individual people is the same.

Nevertheless, a reasonable question is whether we can understand compactness via some simpler geometric approach, analogous to any of the existing measures. The common difficulty of explaining how we as humans (or statistical models that approximate them) perform sophisticated tasks — recognizing a friend's face, developing a scientific hypothesis, judging compactness when we see it, etc. — is known as "Polanyi's paradox," that "we know more than we can tell" (Autor, 2015; Polanyi, 1966). We have studied, in considerable detail, how to simplify our measure and find that indeed the simplest way to know what we see is merely to look or to use our measure. A theoretically simpler version may even be an illusory goal, since humans use such sophisticated combinations of these mathematical simplifications rather than any one. We analyze this point in four ways, and then discuss whether other approaches to this question might be possible.

First, we offer a direct answer from our extensive qualitative analyses of the outputs of our approach along with the features that are most predictive: Our measure of compactness favors districts that are *squarish*, *with minimal arms*, *pockets*, *islands*, *or jagged edges*. (We use "squarish" rather than "circle-like" because many real districts are approximately square-shaped but almost none resemble circles.) Importantly, no existing compactness measure estimates a theoretical quantity that can reasonably be described in this way.

Second, we offer illustrations of the nature of the agreements and disagreements between our measure and each of the seven existing measures we discussed in Section 2. For each existing measure, we construct a 2×2 cross-tabulation of example districts that reflect agreements (compact and noncompact) and disagreements (where the existing measure says noncompact and ours compact, and the reverse). We array horizontally the four cells of this 2×2 table for each measure in a row in Table 2. To generate this table, we define "compact" districts as having a predicted compactness rank in the top 15 (of 100) and "noncompact" as 85 or lower. (If no district appears in a cell of the cross-tabulation, we expand our definition from 15 and 85 to 20 and 80, etc.) Then, to avoid cherry picking, we choose the first in alphabetical order among all districts defined by each cell in each table.⁹

[Table 2 about here.]

The results in Table 2 are striking. The agreements appear in the first two columns: Column one includes seven obviously compact districts, and column two includes seven clearly noncompact districts. The last two columns reflect disagreements between our measure and an existing one. The first of these (in the third column) are districts that our measure indicates are noncompact and an existing measure says are compact. Most human observers agree with our measure (by design) that these are in fact highly noncompact districts. Similarly, the final column includes districts judged as noncompact by an existing measure, but compact by ours. This table clearly reveals how each existing measure picks up important features of the compactness of legislative districts and omits others. The features each measure picks up or misses are those widely discussed in the existing compactness literature as benefits or failures of each measure, since in practice this theoretical literature is using the standard from which our measure was derived (you know it when you see it) to judge their own measures. In contrast, our measure seems to pick up all the features identified throughout the literature as desirable, without obviously missing any feature of a district shape generally seen as important.

Third, do different measures generate different conclusions in practice? The answer here depends on in which legislature the comparison is being made. For any two measures,

⁹We define alphabetical order according to a specific naming convention. All districts receive an identifier which includes state, district set (upper chamber, lower chamber or Congress), district number, and year. For example, Alaska's first congressional district from 2010 is 01_CD_001_2010.

it is easy to draw a districting plan where the measures change the rankings of compactness in *any* arbitrary way. We could also be misled by stacking up data across legislatures — and thus ignoring the bias from heterogeneous treatment effects — in which case we would see that our measure correlates quite low with most measures but at about 0.9 for convex hull and Polsby-Popper, and similarly high correlations for the naive average of all measures. In fact, the only coherent way to answer the question is to use real world legislatures, which is the context in which comparisons matter and, as it turns out, where differences are significant. To pick an extreme case from the current US Congressional map, Polsby-Popper correlates with our measure (i.e., the measure any human observer would choose when evaluating districts) at 0.95 in Indiana's 1970 map but -0.37 in its 1890 map. We thus study this question more systematically by analyzing the 778 legislatures from our data with unique sets of district boundaries (i.e., for every available state, legislative chamber, and year; e.g., Alabama State Senate in 1962). Comparisons across measures in court mostly depend on which district or plan ranks highest and so we compute the percent of times, across data sets, where each existing measure has the highest correlation with our measure. The measure that winds up in the top position most often is Convex Hull, but this occurs in only 54.5% of the data sets — followed by the Polsby-Popper in 31.0%, Grofman in 6.2%, Y-axis Symmetry in 1.9%, Reock and X-axis Symmetry at 1.6% and 1.5%, and Boyce-Clark at 0.6%; even measures such as the area of the minimum bounding circle and the number of discontiguous polygons correlate most highly sometimes. In other words, any existing measure can come out on top in approximating our measure depending on the particular features of the set of district shapes that make up the legislature, and so none of these measures alone can be used as a simpler replacement with our measure of what people know when they see, without checking the relationship first (see Supplementary Appendix E).

Finally, the best practice in choosing predictive models, which we followed, involves finding the most parsimonious model that predicts accurately; as such, we are by definition unable to find an even more parsimonious model without giving up predictive accuracy. Thus, we searched for a more parsimonious model that degraded performance by only a small amount. Unfortunately, we found no large discontinuity in the relationship between parsimony and performance. A straightforward principal component analysis of the existing measures also does not yield a simple solution.

In summary, this section demonstrates that none of the existing measures, and no measure we could find, offer a simple geometric representation for what humans know when they see. To be clear, however, we have not proved that creating such a measure is impossible. We thus leave this as an open question and encourage future researchers to seek such a simplifying geometric definition, if that turns out to be possible.

6 Concluding Remarks

We conclude that the measure derived here reflects the underlying viewpoint held about the concept of compactness by everyone from educated Americans to public officials, judges, and justices. This measure appears to confirm and reflect the single, universally recognizable standard implicit in legal compactness requirements of state constitutions, federal and state legislation, and court decisions. Although "we know more than we can tell" about how humans perceive compactness, this measure quantifies "what we know when we see." The measure is also visibly different (as per Table 2) from any existing measure and, by design, much closer to how human beings perceive compactness.

Approaches developed here for measuring an ill-defined concept that you know only when you see may also be applicable to other difficult-to-define concepts. These include measurement by full ranking rather than paired comparisons, which saves time and turns out, in our application, to have much higher levels of intra- and intercoder reliability; the incorporation in a model rather than replacement of most existing measures and approaches; and formalization into a statistical model of an approach that predicts the views of a wide range of different types of people.

The key aspect of our approach here is defining the concept of interest separately from the measure used to estimate it, so that our measure becomes vulnerable to being proven wrong and, as a result, our approach can improve over time. In this light, we encourage others to take up this challenge and improve on the methods we propose, and develop

22

statistical methods that outperform ours; this may now be possible, as clear performance standards now exist. New features measuring compactness can also be included in our approach as additional covariates in our statistical model, which may well be improved.

We hope the large collection of compactness data we make available with this paper (for 17,896 state legislative and congressional districts) and software that makes it easy to compute compactness on any new district enable future researchers to study a wide range of questions related to this crucial concept (see Appendix E). As well, we hope that having a single measure of compactness that all agree on will begin to constrain some aspects of unbridled advocacy during the redistricting process and subsequent litigation.

Appendix A How Paired Comparisons Fails

The method of paired comparisons has been touted for more than a century and a half for its two key advantages. First, this approach puts fewer demands on survey respondents than asking respondents to do a full ranking. That is, to produce a ranking of n items requires the choice among n! possible rankings, whereas the same information can be elicited with only $\binom{n}{2}$ paired comparisons. This is not trivial since $n! \gg \binom{n}{2}$; for example, with n = 20, we have $20! = 2.4 \times 10^{18}$, or 2 quintillion possible rankings, whereas $\binom{20}{2} = 190$ paired comparisons is large but still manageable in a single survey (and may even be reduced; see Mitliagkas, Gopalan, Caramanis, and Vishwanath 2011). For these reasons, Converse and Presser (1986, p.28) comment on a historical example with only 13 items: "Tasks of this scope were soon seen as much too difficult..., and in our own time, rank orders of this size are all but invisible in the literature". Thus, if full ranking is used, the best practice has been "not to use lists longer than three or four items" (Gideon, 2012).

Second, Thurstone's approach only requires simple questions that are easy to understand, concrete, and specific. With it, we ask a respondent which among a pair of legislative districts is more compact, and then repeat this simple question multiple times with different pairs of districts. Then, after eliciting information in this manner, the researchers combine these binary decisions into a ranked scale (using Guttman scaling or a more sophisticated approach accounting for measurement error; e.g., Mitliagkas, Gopalan, Caramanis, and Vishwanath 2011). The method assumes all respondents will use the same unidimensional scale to make their choices for all their paired comparisons (an issue we return to). The supposed advantage of this approach is that respondents are asked only what they know (a paired comparison) and researchers do what they are better at, which is taking on the complicated task of inferring the underlying full ranking from all the elicited information.

To apply this method, we conducted multiple iterated rounds of pre-testing and cognitive debriefing while adjusting question wording and how the districts appeared¹⁰. But despite dozens of trials over many months, testing numerous variations, and with a wide range of research subjects, online and in person, our inter- and intracoder reliability statistics were rarely much above random chance. To see what we found, consider a simple experiment with 40 respondents (in this case on Amazon's Mechanical Turk), each asked to choose the more compact district from each of twenty pairs, producing a 20-length binary decision vector. This survey enabled us to compare the percent agreement among the 20 decisions for each of $\binom{40}{2} = 780$ pairs of respondents. Figure 7 gives a histogram of these percent agreements (in blue, marked "paired", computed as a density estimate). For comparison, we also generate a placebo test, under the null hypothesis of no agreement, by randomly generating 780 pairs of 20-length vectors and computing from them the percent agreement and plotting its histogram (white with a black outline, marked "Random"). (We discuss the "Ranking" figure in the next section.)

[Figure 7 about here.]

As expected when comparing coin flips, the random placebo percent agreement is centered at 50%. In contrast, the paired comparison percent agreement histogram is shifted farther to the right than the placebo histogram, but the mean only moves to 54%, leaving the two distributions with considerable overlap. Put differently, the best we could do with the method of paired comparisons, even before the step of turning paired decisions into rank orders, is results with unacceptably low levels of intercoder reliability.

¹⁰All districts are visualized at maximally high resolution to ensure that no features such as coastline are lost.

We now rule out the possibility that these results are due to different people having incompatible notions of compactness by studying intracoder reliability. To do this, we waited two weeks, randomly shuffled the order of the 20 paired comparison questions, and administered the survey to the same people. (Of the 40 people, only one mentioned, on post-survey cognitive debriefing, that "some" of the districts may have been the same as the first week.)

These results appear in Figure 8 (also as a blue histogram marked "Paired") and are more distinct from the random placebo test (in white with a black outline marked "Random") than with intercoder reliability in Figure 7, as would be expected. The mean of the paired comparison histogram is now at 65% agreement, although the overlap with the random distribution is still large. (We discuss the third histogram in the next section.)

[Figure 8 about here.]

We thus conclude that these standard, best practice approaches are inadequate, at least for our application, and turn to an alternative. See Section 3.

References

- Altman, Micah (1998): "Modeling the effect of mandatory district compactness on partisan gerrymanders". In: *Political Geography*, no. 8, vol. 17, pp. 989–1012.
- Altman, Micah and Michael P McDonald (2012): "Redistricting principles for the twentyfirst century". In: *Case W. Res. L. Rev.* Vol. 62, p. 1179.
- Angel, Shlomo and Jason Parent (2011): "Non-compactness as voter exchange: Towards a constitutional cure for gerrymandering". In: Northwestern Interdisciplary Law Review, vol. 4, p. 89.
- Ansolabehere, Stephen and Maxwell Palmer (2016): "A Two Hundred-Year Statistical History of the Gerrymander". In: *Ohio St. LJ*, vol. 77, pp. 741–867.
- Ansolabehere, Stephen and James M Snyder Jr (2012): "The effects of redistricting on incumbents". In: *Election Law Journal*, no. 4, vol. 11, pp. 490–502.
- Autor, David (2015): "Polanyi's paradox and the shape of employment growth". In: Federal Reserve Bank of St. Louis: Economic Policy Proceedings, Reevaluating Labor Market Dynamics, pp. 129–177. URL: j.mp/PolanyiP.
- Barabas, Jason and Jennifer Jerit (2004): "Redistricting principles and racial representation". In: *State Politics & Policy Quarterly*, no. 4, vol. 4, pp. 415–435.
- Blasius, Jörg (2012): "Comparing Ranking Techniques in Web Surveys". In: *Field Methods*, no. 4, vol. 24, pp. 382–398.

- Boyce, Ronald R and William AV Clark (1964): "The concept of shape in geography". In: *Geographical Review*, no. 4, vol. 54, pp. 561–572.
- Cain, Bruce (1984): *The Reapportionment Puzzle*. Berkeley: University of California Press.
- Chen, Jowei and Jonathan Rodden (2013): "Unintentional gerrymandering: Political geography and electoral bias in legislatures". In: *Quarterly Journal of Political Science*, no. 3, vol. 8, pp. 239–269.
- Chou, Christine, Steven O Kimbrough, Frederic H Murphy, John Sullivan-Fedock, and C Jason Woodard (2014): "On empirical validation of compactness measures for electoral redistricting and its significance for application of models in the social sciences". In: *Social Science Computer Review*, no. 4, vol. 32, pp. 534–543.
- Converse, Jean M. and Stanley Presser (1986): Survey Questions: Handcrafting the Standardized Questionnaire. Thousand Oaks, CA: Sage Publications.
- David, H. A. (1988): *The Method of Paired Comparisons, 2nd ed.* London: Oxford University Press.
- Fabbris, Luigi (2013): "Measurement scales for scoring or ranking sets of interrelated items". In: *Survey data collection and integration*. Springer, pp. 21–43.
- Fechner, Gustav Theodor (1966): "Elements of psychophysics. Vol. I. [Originally published 1860]". In: *Readings in the history of psychology*.
- Forgette, Richard and Glenn Platt (2005): "Redistricting principles and incumbency protection in the US Congress". In: *Political Geography*, no. 8, vol. 24, pp. 934–951.
- Fryer Jr, Roland G and Richard Holden (2011): "Measuring the compactness of political districting plans". In: *The Journal of Law and Economics*, no. 3, vol. 54, pp. 493–535.
- Gelman, Andrew and Gary King (Sept. 1994): "Enhancing Democracy Through Legislative Redistricting". In: *American Political Science Review*, no. 3, vol. 88, pp. 541–559. URL: j.mp/redenh.
- Gideon, Lior (2012): "The art of question phrasing". In: *Handbook of survey methodology for the social sciences*. Springer, pp. 91–107.
- Grofman, Bernard and Gary King (Jan. 2007): "The Future of Partisan Symmetry as a Judicial Test for Partisan Gerrymandering after LULAC v. Perry". In: *Election Law Journal*, no. 1, vol. 6. http://gking.harvard.edu/files/abs/jp-abs.shtml, pp. 2–35.
- Harris, Curtis C (1964): "A scientific method of districting". In: *Behavioral Science*, no. 3, vol. 9, pp. 219–225.
- Hofeller, Thomas and Bernard Grofman (1990): "Comparing the compactness of California congressional districts under three different plans: 1980, 1982, and 1984". In: *Political Gerrymandering and the Courts*, pp. 281–88.
- Ip, WC, YK Kwan, and LL Chiu (2007): "Modification and simplification of thurstone scaling method, and its demonstration with a crime seriousness assessment". In: *Social indicators research*, no. 3, vol. 82, pp. 433–442.

Kahneman, Daniel (2011): Thinking, fast and slow. Macmillan.

- Kaiser, Henry F (1966): "An objective method for establishing legislative districts". In: *Midwest Journal of Political Science*, no. 2, vol. 10, pp. 200–213.
- Kaufman, Aaron R. (2020): "Implementing novel, flexible, and powerful survey designs in R Shiny". In: *PloS one*, no. 4, vol. 15, e0232424.
- Kaufman, Aaron R., Gary King, and Mayya Komisarchik (2020): Replication Data for: How to Measure Legislative District Compactness If You Only Know it When You See

it. DOI: 10.7910/DVN/FA8FVF. URL: https://doi.org/10.7910/DVN/FA8FVF.

- King, Gary, John Bruce, and Andrew Gelman (1996): "Classifying by Race". In: ed. by Paul E. Peterson. Princeton University Press. Chap. Racial Fairness in Legislative Redistricting. URL: j.mp/Fairrace.
- Krosnick, Jon A. (1999): "Survey Research". In: Annual Review of Psychology, no. 1, vol. 50, pp. 537–567.
- MacEachren, Alan M (1985): "Compactness of geographic shape: Comparison and evaluation of measures". In: *Geografiska Annaler. Series B. Human Geography*, pp. 53– 67.
- Maurer, Daphne, Richard Le Grand, and Catherine J Mondloch (2002): "The many faces of configural processing". In: *Trends in cognitive sciences*, no. 6, vol. 6, pp. 255–260.
- Mitliagkas, Ioannis, Aditya Gopalan, Constantine Caramanis, and Sriram Vishwanath (2011): "User rankings from comparisons: Learning permutations in high dimensions". In: *Communication, Control, and Computing (Allerton), 2011 49th Annual Allerton Conference on.* IEEE, pp. 1143–1150.
- Niemi, Richard G, Bernard Grofman, Carl Carlucci, and Thomas Hofeller (1990): "Measuring compactness and the role of a compactness standard in a test for partisan and racial gerrymandering". In: *The Journal of Politics*, no. 4, vol. 52, pp. 1155–1181.
- Pildes, Richard H and Richard G Niemi (1993): "Expressive Harms, Bizarre Districts, and Voting Rights: Evaluating Election-District Appearances After Shaw v. Reno". In: *Michigan Law Review*, no. 3, vol. 92, pp. 483–587.
- Polanyi, Michael (1966): "The logic of tacit inference". In: *Philosophy*, no. 155, vol. 41, pp. 1–18.
- Polsby, Daniel D and Robert D Popper (1991): "The third criterion: Compactness as a procedural safeguard against partisan gerrymandering". In: *Yale Law & Policy Review*, no. 2, vol. 9, pp. 301–353.
- Prinzmetal, William and Laurie Gettleman (1993): "Vertical-horizontal illusion: One eye is better than two". In: Attention, Perception, & Psychophysics, no. 1, vol. 53, pp. 81– 88.
- Reock, Ernest C (1961): "A note: Measuring compactness as a requirement of legislative apportionment". In: *Midwest Journal of Political Science*, no. 1, vol. 5, pp. 70–74.
- Schwartzberg, Joseph E (1965): "Reapportionment, gerrymanders, and the notion of compactness". In: *Minn. L. Rev.* Vol. 50, p. 443.
- Shi, Jianbo and Carlo Tomasi (1993): *Good features to track*. Tech. rep. Cornell University.
- Stoddart, David R (1965): "The shape of atolls". In: *Marine Geology*, no. 5, vol. 3, pp. 369–383.
- Thurstone, Louis L (1927): "The method of paired comparisons for social values." In: *The Journal of Abnormal and Social Psychology*, no. 4, vol. 21, p. 384.
- Tisdale, Elkanah (1812): "The Gerry-Mander". In: Boston Gazette.
- Wilson, Timothy D and Nancy Brekke (1994): "Mental contamination and mental correction: unwanted influences on judgments and evaluations." In: *Psychological bulletin*, no. 1, vol. 116, p. 117.
- Young, H Peyton (1988): "Measuring the compactness of legislative districts". In: *Legislative Studies Quarterly*, pp. 105–115.

List of Figures

1	Four Districts from the Alabama State House in 2000	29
2	The Underlying Compactness Dimension, from most compact (a) to least	
	compact (d) (all five of the most common compactness measures agree	
	with this ordering). (Districts include, (a) Wyoming State House District	
	42, 2010; (b) Pennsylvania State House District 185, 2010; (c) Oklahoma	
	Congressional District 1, 1950; (d) Louisiana State Senate District 3, 2010.)	30
3	Intercoder Reliability for Full Ranking with 100 districts. Scatterplots are	
	given for the median correlation (top left panel), first quartile (bottom left)	
	and third quartile (top right). A density plot of all correlations, along with	
	a placebo-based density plot appear at the bottom right. Density plots are	
	truncated to reflect the observed support.	31
4	Intracoder Reliability for Full Ranking, following the same heuristics as	
	Figure 3. Density plots are truncated to reflect the observed support	32
5	Cross-Validation of Model Predictions	33
6	Histograms (via density estimates) of correlations between predictions	
	from our model and answers to survey questions from nine different groups	
	of respondents.	34
7	Intercoder Reliability of Thurstone's Paired Comparisons (blue histogram),	
	full ranking (salmon histogram), and a random placebo distribution (white	
	histogram), all using density estimation.	35
8	Intracoder Reliability of Thurstone's Paired Comparisons (blue histogram),	
	full ranking (salmon histogram), and a random placebo distribution (white	
	histogram), all using density estimation.	36



Figure 1: Four Districts from the Alabama State House in 2000.



Figure 2: The Underlying Compactness Dimension, from most compact (a) to least compact (d) (all five of the most common compactness measures agree with this ordering). (Districts include, (a) Wyoming State House District 42, 2010; (b) Pennsylvania State House District 185, 2010; (c) Oklahoma Congressional District 1, 1950; (d) Louisiana State Senate District 3, 2010.)



Figure 3: Intercoder Reliability for Full Ranking with 100 districts. Scatterplots are given for the median correlation (top left panel), first quartile (bottom left) and third quartile (top right). A density plot of all correlations, along with a placebo-based density plot appear at the bottom right. Density plots are truncated to reflect the observed support.



Figure 4: Intracoder Reliability for Full Ranking, following the same heuristics as Figure 3. Density plots are truncated to reflect the observed support.



Figure 5: Cross-Validation of Model Predictions



Figure 6: Histograms (via density estimates) of correlations between predictions from our model and answers to survey questions from nine different groups of respondents.



Figure 7: Intercoder Reliability of Thurstone's Paired Comparisons (blue histogram), full ranking (salmon histogram), and a random placebo distribution (white histogram), all using density estimation.



Figure 8: Intracoder Reliability of Thurstone's Paired Comparisons (blue histogram), full ranking (salmon histogram), and a random placebo distribution (white histogram), all using density estimation.

List of Tables

1	Seven Unique Compactness Rankings of the Same Four Districts: Five	
	Existing and Two New Metrics	38
2	Illustrations of agreements (in the first two columns) and disagreements	
	(in the last two columns) about the degree of compactness between each	
	of seven existing measures and our measure. Each row represents a 2×2	
	table of our measure by an existing measure, with a dichotomized com-	
	pactness summary, displaying one example district in each cell arbitrarily	
	chosen via alphabetical order.	39

	Legislative Districts				
	AL 21	AL 9	AL 62	AL 1	
Convex Hull	1	2	3	4	
Reock	4	3	2	1	
Polsby-Popper	2	3	1	4	
Boyce-Clark	3	4	1	2	
Length/Width	4	2	3	1	
X-axis Symmetry	4	1	2	3	
Significant Corners	3	1	2	4	

Table 1: Seven Unique Compactness Rankings of the Same Four Districts: Five Existing and Two New Metrics



Table 2: Illustrations of agreements (in the first two columns) and disagreements (in the last two columns) about the degree of compactness between each of seven existing measures and our measure. Each row represents a 2×2 table of our measure by an existing measure, with a dichotomized compactness summary, displaying one example district in each cell arbitrarily chosen via alphabetical order.