

CONNECTICUT LAW REVIEW

VOLUME 44

NOVEMBER 2011

NUMBER 1

Article

Rethinking the Indefinite Detention of Sex Offenders

FREDRICK E. VARS

Thousands of sex offenders in the United States are being held indefinitely under civil commitment programs. The analysis in this Article suggests that none (or precious few) belong there. Specifically, in a large dataset, an instrument as good as the one most widely used by experts (the "Static-99") could not identify even one sex offender who met the legal standards for commitment. Supplementing such instruments with additional information does not appear to improve matters, so the failure of the instrument is profoundly disturbing.

There are three possible responses to this failure: (1) improve instruments to meet existing standards; (2) lower the existing standards; or (3) abandon sex offender civil commitment. This Article focuses on the first response, identifying and correcting flaws in the most widely-used instrument. But the greater significance of the Article is to reframe the debate around the other two potential responses. Can we predict the future well enough to justify the indefinite detention of "dangerous" people?

ARTICLE CONTENTS

I. INTRODUCTION	163
II. THREE PROBLEMS WITH THE STATIC-99	167
A. AGE.....	167
B. NORMS	170
C. ERROR AND STANDARDS OF PROOF.....	172
III. A NEW MODEL.....	182
A. DATA	182
B. METHODOLOGY.....	183
C. RESULTS.....	184
1. <i>Age</i>	184
2. <i>Norms</i>	186
3. <i>Error and Standards of Proof</i>	188
IV. DISCUSSION	189
A. LIMITATIONS	189
B. IMPLICATIONS	190
V. CONCLUSION	193
APPENDIX.....	194



Rethinking the Indefinite Detention of Sex Offenders

FREDRICK E. VARS*

[F]rom a legal point of view there is nothing inherently unattainable about a prediction of future criminal conduct.

– Schall v. Martin, 467 U.S. 253, 278 (1984)

Prediction is very difficult, especially about the future.

– Neils Bohr (Danish Physicist)

I. INTRODUCTION

Preventive detention to protect public safety is an old idea, which has recently expanded in scope.¹ Two well-established justifications for civil confinement are mental illness and contagious disease.² A more dramatic example is the now widely discredited internment of Japanese Americans during World War II.³ As the internment case well illustrates, the great difficulty with preventive detention is accurately identifying truly dangerous individuals. This Article focuses on that problem in an area where the assessment of future dangerousness is particularly rigorous: sex offender civil commitment. Even here the prediction problem may be intractable.

Thousands of convicted sex offenders remain in custody after their

* Associate Professor, University of Alabama School of Law (J.D., Yale; A.B., Princeton). I would like to thank Ian Ayres, David Becker, Montré Carodine, Judge Colquitt, Adam Cox, Shahar Dillbary, Ted Donaldson, Caroline Harada, David Kaye, Ron Krotoszynski, Grace Lee, Andy Morriss, Mike Pardo, Hong Min Park, David Patton, Pam Pierson, Hal Stern, Charlie Vars, and Richard Wollert for their help at various stages of this project. Michael Reeves and Jonathan Kolodziej provided excellent research assistance, as did Penny Gibson and the rest of the staff at the University of Alabama law library. Thanks also to Dean Ken Randall and the University of Alabama Law School Foundation for generous research support

¹ Paul H. Robinson, Commentary, *Punishing Dangerousness: Cloaking Preventive Detention as Criminal Justice*, 114 HARV. L. REV. 1429, 1429–31, 1447–49 (2001) (describing the shift over the past several decades from punishment to the prevention of future violations through incarceration).

² *Id.* at 1444 (“Most jurisdictions allow civil commitment of persons who are dangerous because of mental illness, drug dependency, or contagious disease.”).

³ See *Hamdi v. Rumsfeld*, 542 U.S. 507, 542 (2004) (Souter, J., concurring) (“[T]he Emergency Detention Act of 1950 . . . was repealed in 1971 out of fear that it could authorize a repetition of the World War II internment of citizens of Japanese ancestry; Congress meant to preclude another episode like the one described in *Korematsu v. United States*, 323 U.S. 214 (1944).”).

prison terms expire.⁴ As of March 2010, nineteen states and the federal government had laws authorizing the civil commitment of certain sex offenders.⁵ Committed individuals are rarely released.⁶ Moreover, the United States Supreme Court has upheld civil commitment laws against various constitutional challenges.⁷ A common requirement for civil commitment is future dangerousness.⁸ The two primary methods of determining the likelihood of recidivism are clinical judgment and so-called actuarial risk assessment instruments (“ARAI”).⁹ Studies have shown ARAIs to be more accurate.¹⁰ The most widely used ARAI is the Static-99,¹¹ which is the focus of this Article.

⁴ See Monica Davey & Abby Goodnough, *Doubts Rise as States Hold Sex Offenders After Prison*, N.Y. TIMES, Mar. 4, 2007, at A1 (“About 2,700 pedophiles, rapists and other sexual offenders are already being held indefinitely, mostly in special treatment centers, under so-called civil commitment programs . . .”).

⁵ Keith Matheny, *Releases of Sexually Violent Predators Anger Local Areas*, USA TODAY (Mar. 4, 2010), available at http://www.usatoday.com/news/nation/2010-03-03-predator-housing_N.htm (last accessed July 14, 2010); see also 18 U.S.C. § 4248 (2006). Although the federal statute has been held unconstitutional by some courts, see, e.g., *United States v. Wilkinson*, 626 F. Supp. 2d 184 (D. Mass. 2009), the Supreme Court held it to be a constitutional exercise of power under the Necessary and Proper Clause. *United States v. Comstock*, 130 S. Ct. 1949, 1954 (2010). The Court assumed, without deciding, the statute’s constitutionality under other provisions, like the Due Process Clause.

⁶ John Q. La Fond, *The Costs of Enacting a Sexual Predator Law and Recommendations for Keeping Them From Skyrocketing*, in *PROTECTING SOCIETY FROM SEXUALLY DANGEROUS OFFENDERS: LAW, JUSTICE, AND THERAPY* 288 (Bruce J. Winick & John Q. La Fond eds., 2003).

⁷ E.g., *Kansas v. Crane*, 534 U.S. 407, 412 (2002) (upholding, against a due process challenge, state-imposed preventive detention of dangerous sex offenders, but only if there is a lack of control determination); *Kansas v. Hendricks*, 521 U.S. 346, 371 (1997) (holding that the “Kansas Sexually Violent Predator Act comports with due process requirements and neither runs afoul of double jeopardy principles nor constitutes an exercise in impermissible *ex post facto* lawmaking”).

⁸ E.g., *Hendricks*, 521 U.S. at 357 (“Commitment proceedings can be initiated only when a person ‘has been convicted of or charged with a sexually violent offense,’ and ‘suffers from a mental abnormality or personality disorder which makes the person likely to engage in the predatory acts of sexual violence.’”) (quoting KAN. STAT. ANN. § 59-29a02(a) (1994)); see also Robert Prentky et al., *Sexually Violent Predators in the Courtroom: Science on Trial*, 12 PSYCHOL. PUB. POL’Y & L. 357, 357–58 (2006) (examining the most critical problems that occur at the intersection of law and science in the sexually violent predator commitment law context).

⁹ Debra A. Pinals et al., *Violent Risk Assessment*, in *SEX OFFENDERS: IDENTIFICATION, RISK ASSESSMENT, TREATMENT, & LEGAL ISSUES* 70 (Fabian M. Saleh et al. eds., 2009). Some scholars advocate clinically adjusted actuarial assessments. *Id.* at 55. Others strongly disagree with this approach, and continue to advocate traditional clinical judgment. See Stephen D. Hart et al., *Precision of Actuarial Risk Assessment Instruments: Evaluating the ‘Margins of Error’ of Group v. Individual Predictions of Violence*, 190 BRIT. J. PSYCHIATRY s60, s60, s64 (2007) (concluding that ARAI use should be limited to situations that pose low-risk circumstances because of their inadequacy in making accurate predictions about individuals).

¹⁰ Pinals et al., *supra* note 9, at 54; Marcus T. Boccaccini et al., *Field Validity of the Static-99 and MnSOST-R Among Sex Offenders Evaluated for Civil Commitment as Sexually Violent Predators*, 15 PSYCHOL. PUB. POL’Y & L. 278, 278–81 (2009) (“ARAI designed to predict sexual reoffense . . . clearly outperformed unstructured professional judgment . . .”). But see Thomas R. Litwack, *Actuarial Versus Clinical Assessments of Dangerousness*, 7 PSYCHOL. PUB. POL’Y & L. 409, 409–10 (2001) (examining the relative merits of actuarial versus clinical assessments of dangerousness and concluding that it is premature to replace clinical risk assessments with actuarial assessments).

¹¹ Rebecca L. Jackson & Derek T. Hess, *Evaluation for Civil Commitment of Sex Offenders: A Survey of Experts*, 19 SEX ABUSE 425, 434, 448 (2007); Jacqueline Waggoner et al., *A Respecification of Hanson’s Updated Static-99 Experience Table that Controls for the Effects of Age on Sexual*

The Static-99 is a ten-item instrument. The coding form is shown in Table 1:¹²

Table 1. Static-99 Coding Form

Question	Risk Factor	Codes	Score	
1	Young	Aged 25 or older	0	
		Aged 18 – 24.99	1	
2	Ever Lived With	Ever lived with lover for at least two years?		
		Yes	0	
		No	1	
3	Index Non-Sexual Violence-Any Convictions	No	0	
		Yes	1	
4	Prior Non-Sexual Violence-Any Convictions	No	0	
		Yes	1	
5	Prior Sex Offenses	<u>Charges</u>	<u>Convictions</u>	
		None	None	0
		1-2	1	1
		3-5	2-3	2
		6+	4+	3
6	Prior Sentencing Dates (Excluding Index)	3 or less	0	
		4 or more	1	
7	Any Convictions for Non-Contact Sex Offenses	No	0	
		Yes	1	
8	Any Unrelated Victims	No	0	
		Yes	1	
9	Any Stranger Victims	No	0	
		Yes	1	
10	Any Male Victims ¹³	No	0	
		Yes	1	
	Total Score	Add up scores from individual risk factors		

Translating Static-99 Scores into Risk Categories

Score	Label for Risk Category
0, 1	Low
2, 3	Moderate-Low
4, 5	Moderate-High
6 plus	High

Risk categories can be further translated into recidivism rates based on figures from the sample used to develop the instrument:¹⁴

Recidivism Among Young Offenders, 7 LAW PROBABILITY & RISK 305, 305–06 (2008) (“The original version of Static-99 . . . is the most widely used actuarial table for the prediction of sexual recidivism.”).

¹² R. Karl Hanson & David Thornton, *Improving Risk Assessment for Sex Offenders: A Comparison of Three Actuarial Scales*, 24 LAW & HUMAN BEHAV. 119, 133–35 app. I (2000); Andrew Harris et al., *Static-99 Coding Rules Revised—2003* 1, 3–7, 9, 11, available at <http://www.static99.org> (last visited July 1, 2011).

¹³ The Static-99 is designed for male offenders only. The disparate impact of the Static-99 on homosexual and bisexual offenders is beyond the scope of this Article.

¹⁴ Estimates based on Hanson & Thornton, *supra* note 12, at 129 tbl.5.

Table 2. 15-Year Recidivism By Static-99 Risk Category

<i>Risk Category</i>	<i>Sample Size</i>	<i>Sexual</i>	<i>Violent</i>
Low (0, 1)	257	0.09	0.16
Medium-Low (2, 3)	410	0.18	0.32
Medium-High (4, 5)	290	0.37	0.52
High (6+)	129	0.52	0.59
Total (avg. = 3.2)	1086	0.26	0.37

For example, an individual with a score of six or higher on the Static-99 would have a predicted fifteen-year sexual recidivism rate of 52%.

Notably rare in the vast literature examining the Static-99 are studies addressing the fundamental, bottom-line question: can the Static-99 identify individuals who meet the legal standards for commitment?¹⁵ The tests presented in this Article help fill that important gap. The profoundly disturbing answer—given the central role the Static-99 has played in the commitment of thousands of individuals—is essentially no. That answer calls into serious question the entire enterprise of sex offender commitment and, at a minimum, demands the improvement or replacement of the Static-99.

The general approach of this Article is to develop an instrument that predicts recidivism roughly as well as the Static-99 (and is better in other respects), then to ask whether that instrument can identify individuals who qualify for sex offender civil commitment under the existing legal standards. Part II situates the present study within current literature and outlines three problems with the Static-99.¹⁶ Part III reports the results of a new model created and tested in a large dataset.

First, as the creators of the Static-99 have come to recognize, “the

¹⁵ But see Richard Wollert, *Low Base Rates Limit Expert Certainty When Current Actuarials Are Used to Identify Sexually Violent Predators: An Application of Bayes's Theorem*, 12 PSYCHOL. PUB. POL'Y & L. 56 (2006) (“[T]he best available risk-assessment method (i.e., actuarial testing) eventually points to the conclusion that the recidivism rate for each detainee . . . does not meet the commitment standard.”). Two other studies that come closest are: Hart et al., *supra* note 9, at s61–s63 (showing how the approximation uses a social science, not legal, standard), and Eric S. Janus & Paul E. Meehl, *Assessing the Legal Standard for Prediction of Dangerousness in Sex Offender Commitment Proceedings*, 3 PSYCHOL. PUB. POL'Y & L. 33, 40, 60 (1997) (relying on assumptions rather than individual-level data).

¹⁶ There are others. E.g., Melissa Hamilton, *Public Safety, Individual Liberty, and Suspect Science: Future Dangerousness Assessments and Sex Offender Laws*, 82 TEMPLE L. REV. 1, 2–11 (forthcoming 2011), available at <http://ssrn.com/abstract=1580016> (analyzing whether future dangerousness assessments using actuarial tools are responsive to standards contained in sexually violent predators laws). See generally BERNARD E. HARCOURT, *AGAINST PREDICTION: PROFILING, POLICING, AND PUNISHING IN AN ACTUARIAL AGE* 2–3 (2007) (challenging the growing reliance on actuarial methods).

original Static-99 did not sufficiently account for age at release.”¹⁷ The creators have proposed a fix. Part III of this Article outlines a better fix using a fifteen-state dataset: allowing for age effects throughout the range rather than using arbitrary cut-offs.

The second shortcoming, also acknowledged by the originators, is that sexual recidivism rates have fallen since the norms were established. New, lower norms are needed. In Part III, this Article provides one more data point in support of that conclusion and suggests an approach that would more seamlessly adjust to changing crime rates: updating the instrument as soon as new data become available.

Third, and most fundamentally, the Static-99, even as modified, fails to report uncertainty in predicted recidivism rates that is essential to determine whether an individual sex offender meets the commitment threshold according to the applicable standard of proof. Part II explains this problem and classifies each jurisdiction with sex offender civil commitment according to commitment standard and standard of proof. This Article’s alternative prediction model in Part III quantifies the effect of failing to account for prediction error and demonstrates how this shortcoming may (or may not) be overcome. The technical solution is relatively straightforward; the resulting problem, however, is that essentially no one qualifies for commitment. This is the most important finding of this Article: an instrument as good as the Static-99 largely fails to identify any individuals who met the legal standards for commitment. There may be nothing “inherently unattainable” about predicting future behavior, but, in this corner of the real world at least, it is more difficult than present practice admits.

Part IV discusses limitations of the present study—most notably, a short follow-up period—and identifies implications beyond sex offender civil commitment. Actuarial risk assessment is pervasive, particularly in the area of criminal justice.¹⁸ The Static-99 has perhaps more empirical grounding than other widely used instruments, and still it appears to fall short. To illustrate the potentially far-reaching implications of this Article, similarities to one other very popular instrument, the LSI-R, used for parole and other purposes, are highlighted.

II. THREE PROBLEMS WITH THE STATIC-99

A. *Age*

Older people commit fewer crimes. There is a long-recognized inverse

¹⁷ Leslie Helmus et al., *Static-99R: Revised Age Weights* 1, 7 (Oct. 5, 2009), available at <http://www.static99.org> (last visited July 1, 2011).

¹⁸ HARCOURT, *supra* note 16, at 2.

relationship between age and recidivism generally.¹⁹ The review of sexual recidivism, which ultimately led to the Static-99, listed young age as a factor.²⁰ The Static-99 does account for young age: as shown in Table 1, individuals less than twenty-five years old receive an additional point. However, many have criticized the Static-99 for failing to account for age throughout the entire lifespan.²¹ Substantial empirical evidence has accumulated showing that the risk of sexual recidivism declines with age well above twenty-five.²² The general conclusion has been that “recidivism rates decrease in a linear fashion with age-at-release.”²³

The largest field validity test to date found that the Static-99 was not a significant predictor of violent or sexually violent recidivism after controlling for age at release, prior arrests, and release type (mandatory supervision versus discharge).²⁴ Of particular relevance for this Article, age was a highly significant predictor—better than the Static-99.²⁵ On the other hand, when the analysis was limited to sexually violent recidivism, age was not a significant predictor, but the Static-99 score was.²⁶ It should be noted that release type, which was highly significant in both analyses, was apparently determined in part by the Static-99 score.²⁷ Static-99 score also appears to have been used as a screening device elsewhere in the process.²⁸ Using the Static-99 in both these ways would tend to artificially reduce its observed predictive power. Still, the overall finding is that age tends to add predictive power beyond the Static-99.

Hanson and Thornton (with two co-authors) responded to this growing evidence on the Static-99 website.²⁹ First, they confirmed through regression analysis that age was a significant predictor of recidivism even while controlling for the Static-99 score.³⁰ Second, they formulated a new

¹⁹ JOHN MONAHAN, PREDICTING VIOLENT BEHAVIOR: AN ASSESSMENT OF CLINICAL TECHNIQUES 32, 105–07 (1981); see also Robinson, *supra* note 1, at 1451.

²⁰ R. Karl Hanson & Monique T. Bussière, *Predicting Relapse: A Meta-Analysis of Sexual Offender Recidivism Studies*, 66 J. CONSULTING & CLINICAL PSYCHOL. 348, 351 (1998).

²¹ E.g., TERENCE W. CAMPBELL, ASSESSING SEX OFFENDERS: PROBLEMS AND PITFALLS 76 (2004). This criticism applies to all of the five most commonly used ARAIs. See Howard E. Barbaree & Ray Blanchard, *Sexual Deviance Over the Lifespan: Reductions in Deviant Sexual Behavior in the Aging Sex Offender*, in SEXUAL DEVIANCE: THEORY, ASSESSMENT, AND TREATMENT 38 (D. Richard Lewis & William T. O’Donohue eds., 2008).

²² Barbaree & Blanchard, *supra* note 21, at 38, 46, 48–50; LEAM A. CRAIG ET AL., ASSESSING RISK IN SEX OFFENDERS: A PRACTITIONER’S GUIDE 62–67 (2008).

²³ Howard E. Barbaree et al., *The Development of Sexual Aggression Through the Life Span: The Effect of Age on Sexual Arousal and Recidivism Among Sex Offenders*, 989 ANN. N.Y. ACAD. SCI. 59, 67 (2003). Accord Leam A. Craig, *The Effect of Age on Sexual and Violent Recidivism*, 55 INT’L J. OFFENDER THERAPY & COMPARATIVE CRIMINOLOGY 75, 77 (2009).

²⁴ Boccaccini et al., *supra* note 10, at 298 tbl.4.

²⁵ *Id.*

²⁶ *Id.*

²⁷ *Id.* at 292–93, 293 tbl.2.

²⁸ *Id.* at 305.

²⁹ Helmus et al., *supra* note 17, at 1–6.

³⁰ *Id.* at 2–3.

scoring system for age—specifically, ages 18-34.9, +1 point; 35-39.9, 0 points; 40-59.9, -1 point; and 60 or older, -3 points.³¹ Third, they reran the regressions using the modified Static-99 scores. Age was no longer statistically significant,³² leading Hanson and Thornton to conclude that “the original Static-99 did not sufficiently account for age at release, whereas the revised scale does.”³³

This response is unsatisfactory. To be sure, four age categories are better than two, but the real question is why categorize at all? Why not just include age as a continuous variable and let the regression equation assign it the weight that optimizes the model’s predictive power? That is the approach of this Article: to propose a methodology rather than a universal solution. Presumably, the reason Hanson and Thornton resist this approach is attachment to the notion that their instrument needs to be simple enough to be performed with pencil, paper, and no calculator. But Hanson and Thornton have already elsewhere suggested movement away from this model: a computerized coding form, for example, could virtually eliminate logical and arithmetic errors.³⁴ Computerization could just as easily eliminate the conceptual error of applying crude actuarial methods rather than more powerful statistical techniques like logistic regression.

Failing to use age reasonably is arguably unconstitutional.³⁵ Due process dictates that a police officer “may not choose to ignore information that has been offered to him or her.”³⁶ This does not translate into a constitutional duty to investigate, but it does entail a duty not to turn a blind eye to relevant evidence.³⁷ The same principle should apply to adjudicating sex offender civil commitments. “[R]equirements of notice and hearings are of little significance if the decisionmaker ultimately ignores any information before it.”³⁸ The Static-99, both the original and revised versions, effectively throws out relevant information by lumping individuals into broad age categories.

One response to this argument is that the Static-99 is not the only piece of evidence considered, and the Constitution generally constrains the total package, not each constituent part. Decision-makers are free to factor age into the equation, notwithstanding its inclusion in the instrument. That

³¹ *Id.* at 4–5.

³² *Id.* at 5.

³³ *Id.* at 6.

³⁴ R. Karl Hanson et al., *Predicting Recidivism Amongst Sexual Offenders: A Multi-site Study of Static-2002*, 34 *LAW & HUM. BEHAV.* 198, 208 (2010).

³⁵ Others have argued that use of instruments like the Static-99 with older offenders “could be considered to be discriminatory.” See, e.g., Prentky et al., *supra* note 8, at 376.

³⁶ *Kingsland v. City of Miami*, 382 F.3d 1220, 1229 (11th Cir. 2004).

³⁷ *Logsdon v. Hains*, 492 F.3d 334, 341–42 (6th Cir. 2007).

³⁸ Mark Cordes, *Policing Bias and Conflicts of Interest in Zoning Decisionmaking*, 65 *N.D. L. REV.* 161, 217 (1989) (citing Martin H. Redish & Lawrence C. Marshall, *Adjudicatory Independence and the Values of Procedural Due Process*, 95 *YALE L.J.* 455, 476 (1986)).

may be true, but it is almost certainly not the case that decision-makers' informal consideration of age always accurately reflects the true impact of age on recidivism.³⁹ Of course, the government need not wait for a perfect instrument,⁴⁰ but using one that throws away obviously relevant information seems irrational. A second response is that expert testimony based on the Static-99 may not be state action. The nuances of the state action doctrine are outside the scope of this Article; however, there is at least one state where use of the Static-99 is unequivocally state action: Virginia requires its use by statute.⁴¹

B. Norms

The recidivism rates reported in Table 2 were derived from prisoners released from three penal institutions in Canada and one in the United Kingdom.⁴² This should immediately give pause to those who would rely on Table 2 to estimate the likelihood of recidivism for a prisoner in the United States because "the rate of sexual assault in Canada . . . is more than twice that of the United States."⁴³ Moreover, the prisoners in the normative sample were released between the late 1950s and early 1990s. Crime, including sexual offenses, peaked in the early 1990s and has been declining since then.⁴⁴

Given these facts, it should not be surprising that studies have generally found recidivism rates below the Static-99 normative levels.⁴⁵ Hanson, Thornton, and Leslie Helmus, in more recent and diverse samples, found that "sexual recidivism was two-thirds (66%) the rate of the original sample."⁴⁶ There were significant differences among the samples, so, rather than simply adjust the recidivism rates downward based on the

³⁹ See *infra* notes 151–52 and accompanying text.

⁴⁰ Cf. Boccaccini et al., *supra* note 10, at 306 (explaining that governments use instruments that are most effective for their own contexts).

⁴¹ VA. CODE ANN. § 37.2–903 (2006); see also CAL. ANN. PENAL CODE § 290.04(b)(1) (West 2006) (requiring use of Static-99 for registration of sex offenders).

⁴² Hanson & Thornton, *supra* note 12, at 123–24.

⁴³ John A. Fennel, *Punishment by Another Name: The Inherent Overreaching in Sexually Dangerous Person Commitments*, 35 NEW ENG. J. ON CRIM. & CIV. CONFINEMENT 37, 59 (2009).

⁴⁴ Leslie Helmus et al., *Reporting Static-99 in Light of New Research on Recidivism Norms*, 21 THE FORUM 38, 38 (2009).

⁴⁵ Boccaccini et al., *supra* note 10, at 304; see also PATRICK A. LANGAN ET AL., U.S. DEP'T OF JUSTICE, RECIDIVISM OF SEX OFFENDERS RELEASED FROM PRISON IN 1994 1, 1 (2003), available at <http://bjs.ojp.usdoj.gov/content/pub/pdf/rsorp94.pdf> (stating that 5.3% of released sex offenders were rearrested within three years for a sex crime); cf. Shoba Sreenivasan et al., *Predicting the Likelihood of Future Sexual Recidivism: Pilot Study Findings from a California Sex Offender Risk Project and Cross-Validation of the Static-99*, 35 J. AM. ACAD. PSYCHIATRY L. 454, 465 (2007) (explaining that Static-99 under-predicted recidivism at scores of two and three and over-predicted at four through six).

⁴⁶ Helmus et al., *supra* note 44, at 39. But see Grant T. Harris & Marnie E. Rice, *Characterizing the Value of Actuarial Violence Risk Assessments*, 34 CRIM. JUSTICE & BEHAV. 1638, 1643 (2007) (finding no support in five sources cited for the proposition that recidivism rates fell along with overall crime rates).

overall results, the authors provided two estimates for each recidivism type and time period: one lower number for “routine” samples and a higher number for “preselected high risk” samples.⁴⁷ Evaluators are advised to report both the low- and high-end values, then to exercise judgment in opining which sample the individual more closely resembles.⁴⁸ However, the authors concede that the preselection factors that would place an individual into one of the two categories “are not fully known and would vary across samples.”⁴⁹

This “New Norms” article has called into serious doubt use of the Static-99. In *State v. Rosado*,⁵⁰ the sex offender respondent wanted to introduce his Static-99 score of four in the civil commitment proceedings.⁵¹ The court granted the state’s motion *in limine* to exclude the evidence, stating “in view of the recent development of the new norms, and an entirely new and undeveloped methodology for applying those norms, it cannot be said that the new norms of the STATIC-99 (despite its past acceptance) are now sufficiently understood and accepted in the relevant scientific community under *Frye*”⁵²

Supporters of the Static-99 appear caught between a rock and a hard place: they can either fail to account for changes in recidivism rates and generate flawed estimates, or they can adjust their calculations and risk being excluded because the new adjustment is not yet generally accepted in the scientific community. “[T]he development of ARA [actuarial risk assessment], like all good science, is evolutionary.”⁵³ One way for the Static-99 to evolve is suggested by a section of the “New Norms” article that the *Rosado* court did not mention. In it, the authors provide relative risk values for each Static-99 score up to 9+.⁵⁴ Applying those values to the actual recidivism rate for the relevant population would be one reasonable way to generate estimates, at least until there is time for a validation study in that population.⁵⁵

This Article will outline a methodology that nearly automatically adjusts to changes not just in the overall recidivism rate, but also in the

⁴⁷ Helmus et al., *supra* note 44, at 41 tbls.1 & 2. Presumably reflecting the larger sample size, the authors also reported values for each Static-99 score up to 10+. *Id.*

⁴⁸ *Id.* at 40.

⁴⁹ *Id.*

⁵⁰ 889 N.Y.S.2d 369 (2009).

⁵¹ *Id.* at 372–75.

⁵² *Id.* at 416. *Frye v. United States*, 293 F. 1013, 1014 (D.C. Cir. 1923) (stating the standard for the admissibility of scientific evidence in New York).

⁵³ Eric S. Janus & Robert A. Prentky, *Forensic Use of Actuarial Risk Assessment with Sex Offenders: Accuracy, Admissibility and Accountability*, 40 AM. CRIM. L. REV. 1443, 1445 (2003).

⁵⁴ Helmus et al., *supra* note 44, at tbl.3.

⁵⁵ See Calvin M. Langton et al., *Reliability and Validity of the Static-2002 Among Adult Sexual Offenders with Reference to Treatment Status*, 34 CRIM. JUST. & BEHAV. 616, 638 (2007) (“Clearly, likelihood ratios require examination before recidivism probabilities associated with any risk assessment instrument’s scores for one population are assumed to apply to another population.”).

relative predictive power of included variables. Because the proposed methodology is constant, it would arguably not be subject to challenges like the one that succeeded in *Rosado*.

C. Error and Standards of Proof

If a risk scale is to be used in applied contexts, then it is important that the degree of predictive accuracy is sufficient to inform rather than mislead. Critics could suggest, for example, that a correlation in the 0.30 range is insufficient for decision making because it accounts for only 10% of the variance. Even if such an argument was [sic] correct . . . , most decision makers are not particularly concerned about “percent of variance accounted for.” Instead, applied risk decisions typically hinge on whether offenders surpass a specified probability of recidivism (e.g., >50%).⁵⁶

So wrote Karl Hanson and David Thornton in 2000 reporting the results of an early test of the Static-99, which did in fact show correlations around 0.30.⁵⁷ Their statement may be true for certain low-stakes decisions, but it is frighteningly flawed with respect to sex offender civil commitment, as explained below.

In the same article, Hanson and Thornton reported fifteen-year sexual recidivism above their hypothetical 50% threshold for the Static-99 “high” risk category of 52%.⁵⁸ The implication is that commitment would be proper for individuals in that risk category in jurisdictions applying a 50% threshold. That is false. The question is: how sure are we that an individual in this risk category is more likely than not to reoffend? The percentage of variance accounted for, and hence prediction error, is absolutely critical in making that determination. Hanson and Thornton missed the distinction—explained more than twenty years earlier in the same journal by John Monahan and David Wexler—between the standard of commitment and the standard of proof: “one must prove to a given standard [of proof] only that a specified probability threshold [*viz.*, commitment standard] has been crossed”⁵⁹

⁵⁶ Hanson & Thornton, *supra* note 12, at 129–30. A correlation of 1.0 means a perfect positive fit between the predictor and outcome variables; 0 indicates no relationship. The square of the correlation coefficient is the percentage of variance (or spread) of the outcome variable accounted for by the predictor. Thus, as the quoted passage states, a 0.30 correlation coefficient corresponds to roughly 10% of variance accounted for ($0.30^2 = 9\%$).

⁵⁷ *Id.* at 126 tbl.4

⁵⁸ *Id.* at 129 tbl.5, 130; see also *supra* Table 2.

⁵⁹ John Monahan & David B. Wexler, *A Definite Maybe: Proof and Probability in Civil Commitment*, 2 LAW & HUM. BEHAV. 37, 38 (1978); see also M. Neil Browne & Ronda R. Harrison-

If the standard of proof in civil commitment were merely a “preponderance of the evidence” (“POE”)⁶⁰ then, assuming symmetric error, the distinction would not be important.⁶¹ But the United States Supreme Court, in *Addington v. Texas*,⁶² held that due process requires a higher standard for civil commitment. The “clear and convincing evidence” (“CCE”) standard was found to be sufficient, though perhaps not required.⁶³ That standard has been interpreted as proof with greater than 75% confidence.⁶⁴ If the commitment standard is greater than 50%, then the question is whether it is 75% likely that an individual’s risk of recidivism is above that threshold.⁶⁵ The 52% value in the Static-99 recidivism table may or may not be sufficient evidence of that fact, but given the modest correlation with recidivism (0.30), that would seem quite unlikely.

Since 2000, Hanson and Thornton have not been entirely deaf to this criticism.⁶⁶ The Static-2002 is a refinement of the Static-99, and their revised age-specific recidivism risk tables included, for the first time, 95% confidence intervals (“CIs”).⁶⁷ In other contexts, however, they continue to omit critical error estimates. In a leading multi-site study of the Static-2002, Hanson, Thornton, and a co-author once again reported recidivism

Spoerl, *Putting Expert Testimony in its Epistemological Place: What Predictions of Dangerousness in Court Can Teach Us*, 91 MARQ. L. REV. 1119, 1207–10 (2008) (recognizing the significance of the standard of proof). However, Browne & Harrison-Spoerl would simply multiply the commitment threshold by the standard of proof, *id.* at 1209 n.429, which is inappropriate as explained in the text below.

⁶⁰ The POE standard is sometimes described as requiring that a proposition be more likely than not true. In numbers, this translates into a probability greater than 0.5.

⁶¹ D. Mossman & T. Sellke, *Avoiding Errors About ‘Margins of Error,’* 191 BRIT. J. PSYCHIATRY 561 (2007) (explaining how “statistical decision theory” sometimes leads non-lawyer experts in this area to ignore the heightened standard of proof).

⁶² 441 U.S. 418 (1979).

⁶³ *Id.* at 431–33.

⁶⁴ C.M.A. McCauliff, *Burdens of Proof: Degrees of Belief, Quanta of Evidence or Constitutional Guarantees?*, 35 VAND. L. REV. 1293, 1328 tbl.5 (1982) (survey of 170 federal judges reported a mean, median, and mode of 0.75 for the clear and convincing standard); *see also* *United States v. Fatico*, 458 F. Supp. 388, 410 tbl. (E.D.N.Y. 1978) (reporting a range of 0.6 to 0.75 in survey of eight federal district judges); Fredrick E. Vars, *Toward a General Theory of Standards of Proof*, 60 CATH. U. L. REV. 1, 7–11 (2010) (quantifying the effect of the clear and convincing standard). Quantification is resisted by many. *See, e.g.,* Kevin M. Clermont, *Procedure’s Magical Number Three: Psychological Bases for Standards of Decision*, 72 CORNELL L. REV. 1115, 1147–48 (1987) (“[L]awmakers could begin to state the standards in numerical terms. That, however, would not be wise . . .”). That resistance should be somewhat attenuated in this context. For better or worse, commitment thresholds are framed in probabilistic terms, *see infra* Table 3, and most of the evidentiary work is performed by probabilistic or “actuarial” instruments.

⁶⁵ At least this is the way that sex offender civil commitment statutes are in fact structured: with separate commitment and proof standards. That is not, however, the only possible reading of *Addington*. Arguably, setting the commitment standard below 75% violates *Addington*.

⁶⁶ *See, e.g.,* Janus & Prentky, *supra* note 53, at 1471 (pointing out the “absence of information on standard errors”).

⁶⁷ Waggoner et al., *supra* note 11, at 310–11 tbl.3. As explained later in the text, a confidence interval is another, and more useful, measure of the precision of prediction. *See infra* p. 178.

rates for each score, omitting CIs.⁶⁸

Reporting CIs (when they do) is a substantial improvement over the original reports. It falls short, however, of answering the key question. The CIs reported by Hanson and Thornton are group intervals; the legally relevant statistic is the individual interval. In other words, the Static-99 creators are telling us how confident we can be that the recidivism rate (e.g., 52%) accurately reflects the rate for the group of individuals in this risk category (group). A civil commitment decision-maker needs to know how likely it is the individual before it meets the commitment standard (individual).

Stephen Hart and his colleagues have estimated an individual 95% CI on the 52% reported recidivism rate of between 6% and 95%.⁶⁹ In other words, if we had a large sample of individuals in the “High” Static-99 risk category, 95% of them would have a recidivism risk somewhere between 6% and 95%. Some have concluded from this statistic that actuarial risk assessment, and perhaps sex offender commitment generally, should be eliminated.⁷⁰ That may well be the correct conclusion, but the wide confidence interval alone does not decide the issue.

First, and fundamentally, there is arguably a lack of equivalence between the statistical concept of a confidence interval and the legal concept of a standard of proof.⁷¹ As Professor David Kaye observes, “the confidence coefficient is not the probability that [a parameter] lies within the lonely interval we observed. Rather, it is the long run frequency with which various and varied CIs would cover the unknown value for [the parameter].”⁷² Two additional facts arguably are needed to bridge the divide: (1) the probability prior to the subject evidence; and (2) the probability of the evidence under the alternative hypothesis.⁷³ However,

⁶⁸ Hanson et al., *supra* note 34, at 210 tbl.7.

⁶⁹ Hart et al., *supra* note 9, at s62 tbl.2.

⁷⁰ Fennel, *supra* note 43, at 39, 56, 61. Some suggest that the imprecision of ARAIs may render them inadmissible under the standards for expert or scientific evidence. Hart et al., *supra* note 9, at s64. Others contend that ARAIs clear present admissibility hurdles. For an example of this contention, see generally Janus & Prentky, *supra* note 53; see also Christopher Slobogin, *Dangerousness and Expertise Redux*, 56 EMORY L.J. 275, 293–96 (2006) (discussing the court’s “nonchalance” toward prediction testimony). Even if ARAIs are admissible, their imprecision obviously goes to weight and whether dangerousness has been established with the requisite certainty.

⁷¹ See *Turpin v. Merrell Dow Pharmaceuticals, Inc.*, 959 F.2d 1349, 1353 n.1 (6th Cir. 1992) (“The confidence interval is not a ‘burden of proof’ in the legal sense; rather, it is a common sense mechanism upon which statisticians rely to confirm their findings and to lend persuasive power within their profession.”).

⁷² D.H. Kaye, *Apples and Oranges: Confidence Coefficients and the Burden of Persuasion*, 73 CORNELL L. REV. 54, 62 (1987). This connection to frequency is why the approach used to generate CIs is called “frequentist.” *Id.* at 64; see also James M. Curran, *An Introduction to Bayesian Credible Intervals for Sampling Error in DNA Profiles*, 4 LAW PROBABILITY & RISK 115, 116 (2005).

⁷³ These facts are terms in Bayes Theorem. David H. Kaye, *Statistical Significance and the Burden of Persuasion*, LAW & CONTEMP. PROBS., Autumn 1983, at 13, 23. See generally Stanford

the large sample size in this study ($n \approx 9000$) suggests convergence between the frequentist logit CIs presented and the likely results of the alternative, Bayesian methodology.⁷⁴ Thus, although recognizing valid criticisms, this Article equates CIs with standards of proof,⁷⁵ at least as a heuristic device.⁷⁶

Second, as Hart and his fellow authors conceded,⁷⁷ and others further elaborated,⁷⁸ their methodology for estimating CIs had shortcomings. Rehearsing those criticisms here would serve little point, as none of them apply to the CIs calculated in this Article. Hart et al. responded to their critics by explaining that they could have used better methodology if ARAIs were based on logistic regression rather than actuarial methods and if Static-99 data were made publicly available.⁷⁹

Finally, 95% is the conventional spread in social science, but the standard of proof in this context requires a different statistic. Jurisdictions are split between requiring proof beyond a reasonable doubt (“BRD”) and proof by CCE.⁸⁰ As noted above, the latter standard corresponds roughly to 75% certainty. BRD is also generally not quantified by courts, but a

Encyclopedia of Philosophy, Bayes’ Theorem, <http://plato.stanford.edu/entries/bayes-theorem> (last visited Sept. 9, 2011).

⁷⁴ Kaye, *supra* note 72, at 69–70; *see also* M.J. Bayarri & J.O. Berger, *The Interplay of Bayesian and Frequentist Analysis*, 19 STAT. SCI. 58, 71 (2004) (“Bayesian and frequentist asymptotic answers are often (but not always) the same.”); Gauri Sankar Datta et al., *Bayesian Prediction with Approximate Frequentist Validity*, 28 ANNALS OF STAT. 1414, 1414 (2000) (“It is . . . shown that, for any given prior, it may be possible to choose an interval whose Bayesian predictive and frequentist coverage probabilities are asymptotically matched.”).

⁷⁵ For further discussion of the application of confidence intervals to legal contexts see Neil B. Cohen, *Confidence in Probability: Burdens of Persuasion in a World of Imperfect Knowledge*, 60 N.Y.U. L. REV. 385, 399 (1985).

⁷⁶ *See* Neil B. Cohen, *Conceptualizing Proof and Calculating Probabilities: A Response to Professor Kaye*, 73 CORNELL L. REV. 78, 93 (1987) (“I believe that the confidence interval analogy performs well as a heuristic for the decision making process.”); *see also* Cohen, *supra* note 75, at 417 (“[T]o the extent that the model serves a heuristic rather than a technical function, these problems do not detract from the model’s usefulness”). It is worth noting that even Cohen’s most outspoken critic on this point, Professor Kaye, favors the use of interval estimates, just not as equivalents to standards of proof. *See* D.H. Kaye, *Is Proof of Statistical Significance Relevant?*, 61 WASH. L. REV. 1333, 1363–64 (1986) (noting the utility of interval estimation); *see also* Kaye, *supra* note 72, at 68–69 (exemplifying how an equalized test can work like a maximum likelihood test). As an alternative method, in footnote 141 *infra*, I also report results using false positive to false negative error ratios to determine cut-scores.

⁷⁷ Hart et al., *supra* note 9, at s63.

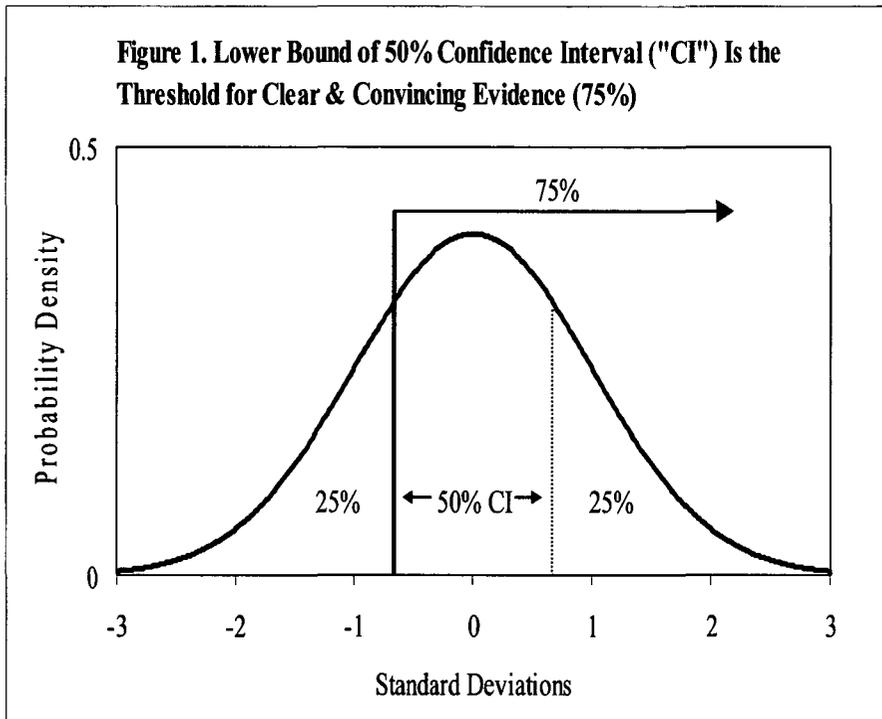
⁷⁸ Harris & Rice, *supra* note 46, at 1648; Grant T. Harris et al., *Shall Evidence-Based Risk Assessment be Abandoned?*, 192 BRIT. J. PSYCHIATRY 154 (2008); Mossman & Selke, *supra* note 61, at 561.

⁷⁹ S.D. Hart et al., *Avoiding Errors About ‘Margins of Error’: Authors’ Reply*, 192 BRIT. J. PSYCHIATRY 561, 561–62 (2007) (correspondence).

⁸⁰ States that demand proof “beyond a reasonable doubt” include Arizona, California, Illinois, Iowa, Kansas, Massachusetts, Missouri, South Carolina, Texas, Washington, and Wisconsin. Nat’l Center for Prosecution of Child Abuse, “Civil Commitment of Sexually Violent Predators,” <http://www.ndaa.org/pdf/Sexually%20Violent%20Predators%20and%20Special%20Sentencing.pdf> (downloaded July 8, 2010). States that require “clear and convincing evidence” include Florida, Minnesota, New Jersey, North Dakota, Virginia. *Id.* For a complete classification of state’s burden of proof requirements, see *infra* Table 3.

survey of 171 judges yielded a mean, median, and mode of 90%.⁸¹

Shifting to a 75% or 90% CI, however, is not the right way to operationalize these standards of proof. Because the commitment question is whether an individual has a risk *greater* than a given threshold, only one tail of the error distribution matters. Assuming symmetric error, applying the CCE standard (75%) therefore requires calculation of the 50% CI: 25% of the error will be below the bottom end of the interval. If that lower bound is above the commitment standard, then commitment is appropriate.⁸² Figure 1 illustrates, assuming a normal error distribution.



This new approach can be applied to the 52% figure with its 6% to 95% range. Assuming a normal error distribution, the lower range of the 50% CI is around 37%. This means there is a 75% chance that the individual recidivism rate is above 37%. That is not enough to clear a “more likely than not” commitment standard, but it is substantially better

⁸¹ McCauliff, *supra* note 64, at 1325, tbl.2. For further information of quantifying BRD, see Vars, *supra* note 64, at 22–23, and Peter Tillers & Jonathan Gottfried, *Case Comment—United States v. Copeland*, 369 F. Supp. 2d 275 (E.D.N.Y. 2005): *A Collateral Attack on the Legal Maxim that Proof Beyond a Reasonable Doubt is Unquantifiable?*, 5 LAW PROBABILITY & RISK 135, 141–51 (2006).

⁸² This approach was suggested by Professor Cohen. See Cohen, *supra* note 75, at 421 (“The difference [between CCE and POE], however, also could be that the clear and convincing standard requires the factfinder to use a higher level of confidence in constructing the interval estimate.”).

than the 6% lower bound of the 95% CI would suggest. For BRD the relevant CI is 80% and the lower bound is around 23%.

Whether values in this range (23–37%) suffice for commitment depends on the commitment standards. Discerning these standards is not always easy. Writing in 1997, Eric Janus and Paul Meehl found that no court or legislature had quantified its standard of commitment.⁸³ They assumed that “highly likely” meant 75% and “likely” meant 50%.⁸⁴ Others more recently claimed that all the state statutes were clear and, despite varying language, set a bar of “roughly 70%.”⁸⁵ Some courts and legislatures have now provided one relatively clear benchmark: Washington, for example, statutorily requires the likelihood of recidivism to be “more probably than not.”⁸⁶ At least one commentator⁸⁷ and two courts⁸⁸ have suggested that “likely” is a lower bar.

Into this morass it is with some trepidation that this Article presents Table 3, illustrating the diversity of both commitment and proof standards:

⁸³ Janus & Meehl, *supra* note 15, at 40, 60. Janus essentially reiterated this point in 2006, calling the thresholds “poorly defined” and “vague.” Prentky et al., *supra* note 8, at 360, 372.

⁸⁴ Janus & Meehl, *supra* note 15, at 41.

⁸⁵ George G. Woodworth & Joseph B. Kadane, *Expert Testimony Supporting Post-Sentence Civil Incarceration of Violent Sexual Offenders*, 3 *LAW PROBABILITY & RISK* 221, 227 (2004).

⁸⁶ WASH. REV. CODE ANN. § 71.09.020(7) (West 2008).

⁸⁷ Jackson & Hess, *supra* note 11, at 439 (“The highest estimate of risk [on the Static-99] is 52% over 15 years. Clearly, 52% barely passes the ‘more likely than not’ criteria established by many states, but may be more than sufficient in a state adopting the language of simply ‘likely.’”); *see also* Thomas Grisso & Paul S. Appelbaum, *Is it Unethical to Offer Predictions of Future Violence?*, 16 *LAW & HUMAN BEHAV.* 621, 627 (1992) (“Given accurate and scientifically supportable predictive testimony about degree of risk, it is up to society (usually its representative on the bench) to determine whether 40%, 30%, or even 20% risk of future violence might reach a threshold justifying a particular legal intervention.”) (internal citations omitted).

⁸⁸ Fennel, *supra* note 43, at 39 (reporting on California and Massachusetts court opinions).

Table 3. Standards of Commitment and Proof by State

<i>Commitment Standard</i>	<i>Proof Standard</i>	
	<i>Clear and Convincing Evidence (75%)</i>	<i>Beyond a Reasonable Doubt (90%)</i>
<i>Likelihood of recidivism</i>		
>50%	Minn., ⁸⁹ N.J. ⁹⁰	Ariz., ⁹¹ Ill. ⁹²
50%	Mo., ⁹³ Neb. ⁹⁴	Iowa, ⁹⁵ Wash., ⁹⁶ Wis. ⁹⁷
<50%	Fed., ⁹⁸ Fla., ⁹⁹	Cal., ¹⁰⁰ Mass. ¹⁰¹
Unspecified	N.H., ¹⁰² N.Y., ¹⁰³ N.D., ¹⁰⁴ Va. ¹⁰⁵	Kan., ¹⁰⁶ S.C., ¹⁰⁷ Tex. ¹⁰⁸

⁸⁹ MINN. STAT. ANN. § 253B.09(1)(A) (West 2007) (“clear and convincing evidence”); MINN. STAT. ANN. § 253B.185 (West 2007) (same); *In re Linehan*, 594 N.W.2d 867, 876 (Minn. 1999) (“highly likely”).

⁹⁰ N.J. STAT. ANN. § 30:4-27.32(a) (West 2008) (“clear and convincing evidence”); *In re Commitment of W.Z.*, 801 A.2d 205, 218 (N.J. 2002) (“highly likely”).

⁹¹ ARIZ. REV. STAT. ANN. § 36-3707(A) (2009) (West) (“beyond a reasonable doubt”); *In re Leon G.*, 26 P.3d 481, 489 (Ariz. 2001) (en banc) (“highly probable”), *vacated on other grounds sub nom.*, *Glick v. Arizona*, 535 U.S. 982 (2002).

⁹² 725 ILL. COMP. STAT. § 207/35(d)(1) (West 2008) (“beyond a reasonable doubt”); *In re Detention of Hayes*, 747 N.E.2d 444, 453 (Ill. App. 2001) (“much more likely than not”).

⁹³ MO. ANN. STAT. § 632.480(5) (West 2006) (“more likely than not”); MO. ANN. STAT. § 632.495(1) (West 2011) (“clear and convincing evidence”).

⁹⁴ NEB. REV. STAT. § 71-1209(1) (2009) (“clear and convincing evidence”); *In re G.H.*, 781 N.W.2d 438, 445 (Neb. 2010) (“more likely than not”).

⁹⁵ IOWA CODE ANN. § 229A.2(4) (West 2006) (“more likely than not”); IOWA CODE ANN. § 229A.7(5)(a) (West 2011) (“beyond a reasonable doubt”).

⁹⁶ WASH. REV. CODE ANN. § 71.09.020(7) (West 2008) (“more probably than not”); WASH. REV. CODE ANN. § 71.09.060(1) (West 2008) (“beyond a reasonable doubt”).

⁹⁷ WIS. STAT. ANN. § 980.01(1m) (West 2007) (“more likely than not”); WIS. STAT. ANN. § 980.05(3)(a) (West 2007) (“beyond a reasonable doubt”).

⁹⁸ 18 U.S.C. § 4247(a)(6) (2006) (“serious difficulty in refraining”); 18 U.S.C. § 4248(d) (2006) (“clear and convincing evidence”); *United States v. Hunt*, 643 F. Supp. 2d 161, 180 (D. Mass 2009) (“[T]his court does not construe the ‘serious difficulty’ criterion for commitment to require proof of any statistical probability of reoffense.”).

⁹⁹ FLA. STAT. ANN. § 394.917(1) (West 2011) (“clear and convincing evidence”); *Hale v. State*, 891 So. 2d 517, 520 (Fla. 2004).

¹⁰⁰ CAL. WELF. & INST. CODE § 6604 (West 2010) (“beyond a reasonable doubt”); *People v. Superior Court (Ghilotti)*, 44 P.3d 949, 968 (Cal. 2002) (stating that “‘likely . . . does not mean the risk of reoffense must be higher than 50 percent,’” but instead means the person “presents a *substantial danger*—that is, a *serious and well-founded risk*—of reoffending”).

¹⁰¹ MASS. GEN. LAWS ch. 123A, § 14(d) (2003) (“beyond a reasonable doubt”); *Commonwealth v. Boucher*, 780 N.E.2d 47, 53 (Mass. 2002) (defining “‘likely’” not as “‘more likely than not,’” but rather as “‘would reasonably be expected’”).

¹⁰² N.H. REV. STAT. ANN. § 135-E:2(VI) (2010) (“potentially serious likelihood”); N.H. REV. STAT. ANN. § 135-E:11(I) (2010) (“clear and convincing evidence”); *State v. Paradis*, 455 A.2d 1070, 1072 (N.H. 1983) (“dangerous”).

¹⁰³ N.Y. MENTAL HYG. LAW § 10.03(e) (McKinney 2011) (“likely to be a danger to others”); N.Y. MENTAL HYG. LAW § 10.07(d) (McKinney 2011) (“clear and convincing evidence”).

Table 3 and the accompanying notes illustrate two important points: (1) there is great diversity *across* states in approaches on both standard of commitment and standard of proof; and (2) with the exception of the relatively precise “more likely than not” standard, the vague statements of commitment standards strongly suggest that there is no uniformity *within* most states either.¹⁰⁹ On the second point, commentators have recommended quantification “so that the true distribution of the risk of error in prediction can be seen.”¹¹⁰ This would reveal the otherwise hidden policy tradeoffs.¹¹¹ The debate about quantification is, as in other contexts, partly about who one wants to make these tradeoffs.¹¹² By adopting numerical standards, legislatures and appellate courts can shift discretion away from trial courts, juries, and testifying experts. Doing so advances the goals of transparency and consistency, but at the price of case-specific flexibility.¹¹³

One argument for flexible standards in this context deserves special attention. The likelihood of recidivism is only one component in determining the gravity of the threat posed by a particular sex offender; also relevant are the magnitude of future harms, their frequency, and their imminence.¹¹⁴ Sex offender commitment statutes generally do not capture these other elements, but case-specific adjustment of the required probability level might (e.g., “menace”). Current ARAIs may aggravate the problem: “existing actuarial methods are optimized to predict the most common but least severe sexual offenses.”¹¹⁵ Of course, ARAIs and statutes could be adjusted to address the problem without sacrificing

¹⁰⁴ N.D. CENT. CODE § 25-03.3-13 (2009) (“clear and convincing evidence”); *In re B.V.*, 708 N.W.2d 877, 882 (N.D. 2006) (stating that defining “likely” as “of such a degree as to pose a threat to others . . . ‘prevents a contest over percentage points and the results of other actuarial tools’”).

¹⁰⁵ VA. CODE ANN. § 37.2-908(C) (2005) (“clear and convincing evidence”); *Shivae v. Commonwealth*, 613 S.E.2d 570, 577 (Va. 2005) (“a menace to the health and safety of others”) (internal quotations omitted).

¹⁰⁶ KAN. STAT. ANN. § 59-29a02(c) (2005) (“menace”); KAN. STAT. ANN. § 59-29a07(a) (2005) (“beyond a reasonable doubt”).

¹⁰⁷ S.C. CODE ANN. § 44-48-30(9) (2010) (“pose a menace”); S.C. CODE ANN. § 44-48-100(A) (2010) (“beyond a reasonable doubt”).

¹⁰⁸ TEX. HEALTH & SAFETY CODE ANN. § 841.003(a)(2) (West 2010) (“likely”); *Beasley v. Molett*, 95 S.W.3d 590, 600 (Tex. App. 2002) (“The term ‘likely,’ as ordinarily defined, means ‘probable.’ Something that is probable is beyond a mere possibility or potential for harm.”). TEX. HEALTH & SAFETY CODE ANN. § 841.062(a) (West 2010) (“beyond a reasonable doubt”).

¹⁰⁹ See Jason A. Cantone, *Rational Enough to Punish, but too Irrational to Release: The Integrity of Sex Offender Civil Commitment*, 57 *DRAKE L. REV.* 693, 712–13 (2009) (noting that vague legal standards may lead to lack of uniformity).

¹¹⁰ Janus & Meehl, *supra* note 15, at 34.

¹¹¹ *Id.*; cf. Grisso & Appelbaum, *supra* note 87, at 627.

¹¹² Vars, *supra* note 64, at 21–22.

¹¹³ *Id.*

¹¹⁴ Janus & Prentky, *supra* note 53, at 1449.

¹¹⁵ Sreenivasan et al., *supra* note 45, at 466.

transparency and consistency, but decision-maker discretion with non-quantified risk thresholds is perhaps a next best solution.

The diversity of approaches shown in Table 3 underscores the importance of the fact that ARAIs like the Static-99 are meaningful only when confidence intervals are properly understood and reported. There are at least six possible combinations of standard of proof and standard of commitment. The state-specific combination must be factored into any ultimate opinion based on an ARAI. And experts are making such judgments as a matter of routine. Ninety-five percent of evaluators reported using the Static-99 “always or most of the time,” and the same percentage “reported that it was either essential or recommended for an evaluator to state an ultimate opinion regarding whether a sex offender meets civil commitment criteria in their final report.”¹¹⁶

Whether or not one favors quantification of the commitment and proof standards, one should favor quantification of the likely error associated with the Static-99 or other ARAIs. The recidivism tables and risk categories are at best misleading in the absence of confidence intervals. No one—not an expert, a trial court, a jury, an appellate court, nor a legislature—can balance the costs and benefits of sex offender commitment without some sense of the error of actuarial prediction.¹¹⁷ Experts are ethically bound to report the limitations of actuarial results.¹¹⁸ Actuarial and other statistical methods have the potential to generate both risk estimates and confidence intervals on those estimates.¹¹⁹ Parts III and IV of this Article realize that potential for a new model and consider what the results mean for risk assessment generally.

This Article follows in the footsteps of Janus and Meehl.¹²⁰ Given certain assumptions about the meaning of commitment standards (e.g., “likely” = 50%; “highly likely” = 75%),¹²¹ accuracy of prediction (0.75),¹²²

¹¹⁶ Jackson & Hess, *supra* note 11, at 434, 435.

¹¹⁷ *Id.* at 439.

¹¹⁸ Grisso & Appelbaum, *supra* note 87, at 630 (explaining that an expert has an ethical duty of “presenting reliable testimony and clearly explaining its limitations”); Stephen D. Hart et al., *A Note on Portraying the Accuracy of Violence Predictions*, 17 LAW & HUM. BEHAV. 695, 696 (1993) (“In the context of psycholegal assessments, unwillingness to qualify one’s confidence in violence predictions or failure to make probabilistic statements regarding the likelihood of future violence is, at best, poor practice; at worst, it is simply unethical.”); Randy K. Otto & John Petrila, *Admissibility of Testimony Based on Actuarial Scales in Sex Offender Commitments: A Reply to Doren*, 3 SEX OFFENDER L. REP. 1, 15 (Dec./Jan. 2002) (“[T]here is an obligation on the part of experts to be as precise as possible not only about their testimony, but about the limitations on the tests that underlie their testimony.”).

¹¹⁹ Janus & Prentky, *supra* note 53, at 1493. Some have suggested that this is impossible. See Gina M. Vincent et al., *The Use of Actuarial Risk Assessment Instruments in Sex Offenders*, in SEX OFFENDERS: IDENTIFICATION, RISK ASSESSMENT, TREATMENT, AND LEGAL ISSUES 71 (Fabian M. Saleh et al. eds., 2009) (asserting that estimates “cannot be done with known precision”).

¹²⁰ Janus & Meehl, *supra* note 15.

¹²¹ *Id.* at 41.

¹²² *Id.* at 49.

and the base rate of recidivism (20-45%),¹²³ Janus and Meehl concluded that it was possible to meet the commitment standard. Notably, the standard of proof was simply folded into standard of commitment with no downward adjustment in the likelihood of recidivism.¹²⁴ This Article corrects that methodological error¹²⁵ and replicates, for a particular prediction model and dataset, what Janus and Meehl attempted as a matter of theory. This Article makes no assumptions about accuracy or base rate,¹²⁶ but rather lets the data set those values. Nor does this Article assume a single cut-score.¹²⁷ Finally, this Article goes farther than Janus and Meehl by pointing the direction toward better ARAIs.

George Woodworth and Joseph Kadane also examined a particular ARAI—in their case another pre-existing instrument—the Mn-SOST-R.¹²⁸ That instrument shares the deficiencies of the Static-99 as described above. Furthermore, Woodworth and Kadane collapsed standards of commitment into a single, unsupported cut-off and ignored standards of proof entirely.¹²⁹ Despite these limitations, the present Article agrees with and implements their suggestion that logistic regression can improve upon less sophisticated, “actuarial” methods.¹³⁰

In perhaps the closest precursor to the present Article, Richard Wollert applied Bayesian techniques to evaluate several ARAIs, including the Static-99.¹³¹ Wollert apparently followed Janus and Meehl in assuming a recidivism threshold of between 50% and 75%.¹³² He not only found that the studied instruments failed to identify even one individual qualified for commitment, but concluded that the instruments would always fail unless base rate recidivism rose above 25%.¹³³ Wollert again used a single cut-off, but this Article confirms his basic results using corrected commitment criteria and much different methodology.

¹²³ *Id.* at 51.

¹²⁴ *Id.* at 43.

¹²⁵ I use the term “error” descriptively, not normatively. Jurisdictions in fact decouple the commitment standard and standard of proof. Whether that bifurcated approach is defensible (or constitutional) is outside the scope of this Article.

¹²⁶ See Dennis M. Doren & Douglas L. Epperson, *Great Analysis, but Problematic Assumptions: A Critique of Janus and Meehl (1997)*, 13 *SEXUAL ABUSE: J. RES. & TREATMENT* 45, 46–48 (2001) (arguing that the assumed base rate was too low).

¹²⁷ See *id.* at 49–51 (noting that Janus and Meehl assume a single, definitive cut-score).

¹²⁸ See Woodworth & Kadane, *supra* note 85, at 231.

¹²⁹ See *id.* at 227 (“[W]e believe that the legislatures intend and the courts require a probability of roughly 70%.”); *id.* at 239 (treating recidivism percentage equal to or greater than commitment standard as sufficient to justify commitment).

¹³⁰ *Id.* at 239.

¹³¹ Wollert, *supra* note 15.

¹³² *Id.* at 58.

¹³³ *Id.* at 75, 79.

III. A NEW MODEL

A. *Data*

The data for this study are taken from a United States Department of Justice, Bureau of Justice Statistics (“BJS”) database.¹³⁴ That database includes all prior criminal history information and recidivism over a three-year follow-up period for 38,624 sampled prisoners released from prisons in fifteen states in 1994.¹³⁵ The data were chosen for several reasons: (1) the dataset is very large and therefore comes closest to representing the United States as a whole, and (2) it covers a time period in which sex offender commitment was not yet prevalent. Thus, it includes every sex offender, not just those deemed safe enough to release.¹³⁶ All violent sex offenders were included in the BJS study; non-violent sex offenders were sampled. Because the present study includes both types of sex offenders, probability weights were used to adjust for sampling.¹³⁷

The present study is limited to the 10,400 men in the BJS study who were incarcerated for a sex offense immediately prior to their release in 1994. Table 4 sets forth some of their relevant characteristics.

¹³⁴ *Recidivism of Prisoners Released in 1994: [United States]* (ICPSR Study No. 3355).

¹³⁵ U.S. DEP’T OF JUSTICE, BUREAU OF JUSTICE STATISTICS, *RECIDIVISM OF PRISONERS RELEASED IN 1994: [UNITED STATES], CODEBOOK iii* (2002) [hereinafter *CODEBOOK*].

The states are: Arizona, California, Delaware, Florida, Illinois, Maryland, Michigan, Minnesota, New Jersey, New York, North Carolina, Ohio, Oregon, Texas, and Virginia. *Id.* at iv.

¹³⁶ *Id.* at 3. No state in the study had sex offender commitment in 1994, except for Minnesota late in that year. Since they had expressly not been selected for commitment, individuals released in Minnesota after the effective date of the sexual offender commitment law were omitted. See MINN. STAT. ANN. § 253B.185 (West 2007) (effective Sept. 1, 1994).

¹³⁷ See *CODEBOOK*, *supra* note 135, at 12–13. All analyses were rerun without weights. There were only trivial changes in results. This was not unexpected as in most cases fewer than fifty individuals had probability weights not equal to one. Along the same lines, excluding non-violent sex offenders had no significant impact on the results. Such offenders made up about 3% of the total and were mostly serving time for statutory rape or incest. See *infra* Table 4.

Table 4. Summary Statistics (Unweighted)

	<i>Number</i>	<i>Percentage</i>
<u>Offense</u>		
Rape	2407	23.1%
Statutory Rape	282	2.7%
Incest	59	0.6%
Sexual Abuse	3856	37.1%
Child Molestation	3278	31.5%
Sodomy	518	5.0%
<u>Race/Ethnicity</u>		
Black	3238	31.1%
Hispanic	1697	16.3%
<u>Age*</u>		
<25	1340	12.9%
25–35	3733	35.9%
35–50	4141	39.9%
>50	1177	11.3%

*Age is missing for nine individuals.

Due to missing data, roughly 15% of the total sample was excluded: the key regression presented below in Table 5 is based on 8881 instead of 10,400 observations.¹³⁸

B. Methodology

This Article employs the BJS data to shed light on current sex offender commitment practice. This is *not* a direct test of the efficacy of the Static-99, because data including Static-99 scores have not been made publicly available. Rather, more sophisticated statistical tools are used to demonstrate the points made above regarding age, norms, and error. Most fundamentally, the data are used to estimate how many individuals met the legal standards for commitment, and how error of prediction affects that estimate. The primary data analysis tool is logistic regression. It is a commonly-used model in social science for true dichotomous outcomes,

¹³⁸ By far the largest source of missing data is the lack of arrest data. Every convict must have been arrested at least once, so observations without any arrests were dropped. In contrast, I retained observations with one or more arrest charge and additional charges coded as “missing” or “unknown.” A strict reading of the Codebook would exclude such individuals because known negatives should have been coded “not applicable.” See CODEBOOK, *supra* note 135, at 13–14. Such a reading, however, would in almost every case contradict the recorded number of charges per arrest (e.g., variable name = A001CNT).

like recidivism. Technical details of the methodology are presented in the Appendix.

C. Results

1. Age

As described above, the original Static-99 converts the continuous variable age into a dichotomous variable with “Young” equal to one if the offender is less than twenty-five years old at the time of release.¹³⁹ Even the creators of the Static-99 admit that this was a mistake. One obvious question is whether a more refined treatment of age, on its own, can predict recidivism as well as or better than the Static-99. The answer is mixed. Age alone does as well as the Static-99 in predicting violent (both sexual and non-sexual) recidivism, but not as well in predicting sexual (both violent and non-violent) recidivism.

The present study includes two continuous age variables: age at first arrest and age at release. Squared and cubed versions of each are also included to allow non-linear effects.¹⁴⁰ Two logistic regression (or “logit”) models estimated the likelihood of recidivism using these six age variables alone. The mean predicted likelihood for the group that was arrested for a subsequent violent offense was significantly higher than the group that did not recidivate violently: 30.6% versus 24.6% (Cohen’s $d = 0.58$ [95% CI = 0.53, 0.63]). This effect size—a measure of the strength of association between two variables—matches that of the Static-99. A recent meta-analysis of thirty-five studies of the Static-99’s predictions of violent recidivism found a mean Cohen’s d equal to 0.57, with a 95% CI of 0.52 to 0.62.¹⁴¹

Sexual recidivism is a different story. The same Static-99 meta-analysis found a mean Cohen’s d of 0.67 in predicting sexual recidivism with a 95% confidence interval of 0.62 to 0.72. For the age-only logit model described above, the Cohen’s d was only 0.30 (95% CI = 0.23, 0.38). Hence, the Static-99 was significantly better than age alone at predicting sexual recidivism.

But that is not really a fair comparison. The Static-99 includes nine variables other than age. Five variables assess prior involvement with the criminal justice system. The question, then, is can the combination of the two continuous age variables, two criminal history variables, and superior methodology (logistic regression) compete with the Static-99 in predicting

¹³⁹ See *supra* Table 1.

¹⁴⁰ This is in some ways a less constrained version of Prentky et al., *supra* note 8, at 376–77, who concluded that recidivism estimates should be reduced by two percent for every year after age forty. The results reported in Table 5 suggest a large negative effect of age throughout the range of released individuals.

¹⁴¹ See *infra*, Appendix (defining Cohen’s d).

sexual recidivism? The new model adds two sets of variables closely mimicking two items on the Static-99: (1) a dummy variable equal to one if the individual has a prior conviction for a violent offense (along with another dummy equal to one if offense code was missing); and (2) a set of dummy variables based on the number of prior arrests for sexual offenses (0 or 1; 2 or 3; 4, 5, or 6; and 7 or more). Table 5 reports the results.

Table 5. Logit Regression Predicting Rearrest for Sexual Offense

Log likelihood = -2591.4108						
Number of obs = 8881						
Pseudo R2 = 0.0557						
Area under ROC curve* = 0.6605						
Cohen's <i>d</i> = 0.761						
<i>Variable</i>	<i>Coefficient</i>	<i>Std. Error</i>	<i>z</i>	<i>P> z </i>	<i>[95% Conf. Interval]</i>	
Age at Release						
Release	-0.000384	0.000351	-1.09	0.275	-0.00107	0.000305
...squared	3.31E-08	2.53E-08	1.31	0.191	-1.7E-08	8.26E-08
...cubed	-9.33E-13	5.89E-13	-1.58	0.113	-2.1E-12	2.21E-13
Age at First Arrest						
Arrest	1.06E-04	1.52E-04	0.70	0.485	-1.92E-04	4.05E-04
...squared	-2.19E-08	1.36E-08	-1.61	0.108	-4.86E-08	4.81E-09
...cubed	7.86E-13	3.91E-13	2.01	0.045	1.92E-14	1.55E-12
Violent Prior Conviction?						
...yes	-0.220320	0.110261	-2.00	0.046	-0.436428	-0.004211
...missing	0.420117	0.201135	2.09	0.037	0.025900	0.814334
Sexual Prior Arrests						
...0 or 1	-1.494985	0.230609	-6.48	0.000	-1.946970	-1.043000
...2 or 3	-0.875472	0.227590	-3.85	0.000	-1.321540	-0.429404
...4, 5, or 6	-0.289424	0.270223	-1.07	0.284	-0.819052	0.240205
...6 or more	0.023862	0.561787	0.04	0.966	-1.077220	1.124945
Constant	0.719014	1.414187	0.51	0.611	-2.052742	3.490770

*Based on unweighted regression

Perhaps not surprisingly, the most significant predictor of sexual recidivism was a significant number of prior arrests for sexual offenses. But more important for the present Article is the result that better methodology and more sensitive treatment of age more than compensated for omitting seven of the ten Static-99 items: the Cohen's *d* for this expanded model was 0.76, which is above the 95% CI reported by the Static-99 meta-analysis. To be sure, the CIs overlap substantially, but it appears fairly safe to say that the new model does as well or better than the Static-99.

However, the analysis to this point has been unfair in at least one way to the Static-99. The meta-analysis reviewed applications of the Static-99 to populations other than the ones with which the instrument was

developed. In contrast, the logit models have been evaluated with the same data that was used to construct them. This out-of-sample versus in-sample comparison is not apples to apples. A properly specified model will generally perform best in the construction sample. To correct this bias, the data were split into two parts. The model was constructed using data from every state other than California and its predictions were tested in California. California was chosen because it was home to the largest number of released prisoners, roughly one-third of the total.

The results with respect to violent recidivism stand: the logit model with two sets of age variables does as well as or better than the Static-99 in predicting violent recidivism in new samples. Specifically, the Cohen's *d* for the out-of-sample logit model was 0.56 with a 95% confidence interval of 0.48 to 0.64. This Cohen's *d* is nearly identical to the Static-99.

In contrast, the Static-99 has an edge over the logit model for out-of-sample prediction of sexual recidivism. The Cohen's *d* for logit is 0.53 (95% CI = 0.41, 0.65), as compared with 0.67 (95% CI = 0.62, 0.72) for Static-99. Note that the confidence intervals overlap, so one cannot reject the hypothesis of equivalence. Recidivism rates vary considerably by race and ethnicity in this dataset,¹⁴² so including these variables might increase predictive power, particularly since California's demographics may be unusual. However, after controlling for age and criminal history, race and ethnicity actually *reduce* effect size, strongly suggesting that these characteristics do not belong in the model.¹⁴³

To summarize, the data support using age as a continuous variable with a technique like logistic regression. A logit model on age at first arrest and age at release, together with squared and cubed terms allowing for non-linearity, was better at predicting violent recidivism than the ten-item Static-99. Adding two sets of criminal history variables made the logit model nearly as good as the Static-99 in predicting sexual recidivism. In other words, more sophisticated use of age eliminated the need to collect additional criminal history information or to code the three Static-99 items based on victim characteristics.¹⁴⁴ Of course, even greater predictive power would likely be achievable by including those items in a regression model.

2. Norms

Recall that recidivism rates vary significantly across time and jurisdiction. The present study is no exception. As shown in Table 6, there is wide disparity in the recidivism rates of sex offenders among the

¹⁴² CODEBOOK, *supra* note 135, at 18 tbl.12 (2003).

¹⁴³ Consistent with this conclusion, the coefficient on the black variable failed to achieve statistical significance. However, the Hispanic coefficient was negative and significant ($p = 0.008$).

¹⁴⁴ See *supra* Table 1.

fifteen states included in this study.

Table 6. Rearrest Rate by State and Offense Type

<i>State</i>	<i>Sexual</i>	<i>Violent</i>
Arizona	5.9%	22.8%
California	8.4%	24.6%
Delaware	10.4%	42.9%
Florida	9.7%	28.5%
Illinois	12.2%	39.6%
Maryland	13.5%	40.0%
Michigan	5.7%	13.6%
Minnesota	18.1%	27.1%
New Jersey	6.6%	23.7%
New York	10.0%	26.2%
North Carolina	7.3%	26.6%
Ohio	13.0%	31.2%
Oregon	8.6%	23.8%
Texas	6.7%	23.0%
Virginia	9.3%	27.5%

This disparity suggests that state-specific norms are needed to evaluate an individual's risk of recidivism. Better evidence would be to find significant state effects after controlling for age and criminal history.

Dummy variables for fourteen of the fifteen states (the last, Virginia, was omitted) were added to the sexual recidivism model outlined in the previous section.¹⁴⁵ The coefficients on two states, Maryland and Texas, were negative and statistically significant ($p < 0.05$).¹⁴⁶ A chi-square test rejected the null hypothesis that all the state dummy coefficients were equal ($\chi^2(14) = 54.01$; $p < 0.0001$). The analysis was repeated for violent recidivism. Here again, the independent variables were based on age and criminal history, along with state dummy variables. Five state coefficients were statistically significant at the 5% level (Illinois, Maryland, Michigan, North Carolina, and Texas). As a result, it is not surprising that overall there is a highly significant difference among states ($\chi^2(14) = 106.73$; $p <$

¹⁴⁵ For technical reasons, there must always be an omitted reference category for regression models to work. See ROBERT S. PINDYCK & DANIEL L. RUBINFELD, *ECONOMETRIC MODELS AND ECONOMIC FORECASTS* 106 (3d ed. 1991). See generally *Dummy Variable (Statistics)*, WIKIPEDIA (July 4, 2011, 10:38 AM), [http://en.wikipedia.org/wiki/Dummy_variable_\(statistics\)](http://en.wikipedia.org/wiki/Dummy_variable_(statistics)) (footnote omitted) ("If dummy variables for all categories were included, their sum would equal 1 for all observations, which is identical to and hence perfectly correlated with the vector-of-ones variable whose coefficient is the constant term; if the vector-of-ones variable were also present, this would result in perfect multicollinearity, so that the matrix inversion in the estimation algorithm would be impossible. This is referred to as the dummy variable trap.")

¹⁴⁶ North Carolina had a negative coefficient that came very close to statistical significance ($p = 0.054$).

0.0001).

Many factors could explain the significant state effects. The important point is that significant differences among states persist even after controlling for factors like those included in the Static-99. Hence, those creating risk assessment instruments and those using them should seriously consider state-specific norms.

3. *Error and Standards of Proof*

Again, the logit model described above predicted sexual recidivism within sample as well as, or better than, the Static-99. One great advantage of the logit model is that it is possible to measure the standard error associated with individual predictions. As a result, one can actually test whether each released sex offender exceeded the commitment standard with the requisite degree of confidence. Table 7 summarizes the results.

Table 7. Individuals Who Met the Dangerousness Standards for Commitment (Out of 8881)*

<i>Commitment Standard</i>	<i>Proof Standard</i>	
	<i>Clear and Convincing Evidence (75%)</i>	<i>Beyond a Reasonable Doubt (90%)</i>
<i>Likelihood of Recidivism</i>		
> 75%	0	0
> 50%	0	0
> 25%	217 2.4%	201 2.3%

*Using predictions and standard errors of model summarized in Table 5

The most striking result is that *not one* of the 8881 released sex offenders was more likely than not to be rearrested for a sexual offense—even at the lower CCE standard. This means that, using the instrument alone, no one met the dangerousness threshold used in half of the jurisdictions with sex offender commitment.¹⁴⁷

However, at least four jurisdictions—California, Florida, Massachusetts, and the federal government—set the bar lower than that: a less than 50% chance of recidivism (which, again, was arbitrarily set at

¹⁴⁷ See *supra* Table 3. Some have argued that violent recidivism is a better measure among sex offenders of the conduct civil commitment is designed to prevent. See Marnie E. Rice et al., *Violent Sex Offenses: How Are They Best Measured from Official Records?*, 30 LAW & HUM. BEHAV. 525, 536 (2006). The analysis was repeated using the logit model described above predicting violent recidivism using six age variables. Despite a much higher base rate of violent recidivism (about 26%), no one qualified for commitment at the 75% threshold and around 0.2% at the 50% threshold (out of 9,015 individuals, twenty-three at CCE and nineteen at BRD). However, solid majorities cleared the 25% hurdle.

25%). About 2.3% of individuals (201 of 8881) met this standard beyond a reasonable doubt. Among this most dangerous group, the actual recidivism rate was very close to 40%. In other words, if these individuals had been committed, three people who would not have reoffended would have been detained for every two recidivists. This analysis also showed that almost 90% of recidivists still would have been released.¹⁴⁸

The contribution of error can be quantified. How many individuals would have qualified for commitment if the error associated with predictions were ignored, as the creators of the Static-99 originally advocated?¹⁴⁹ Still none met the 50% and 75% commitment standards. At the lowest threshold (25%), however, 242 individuals would have qualified. In other words, properly factoring in error and applying a higher standard of proof can reduce commitments identified by the logit model by up to 17% (242 to 201). The error of individual prediction associated with an actuarial tool like the Static-99 is likely greater because it rounds variable effects and lumps individuals into rough categories.¹⁵⁰

IV. DISCUSSION

A. *Limitations*

This is not a direct test of the Static-99. Nonetheless, the findings described above illuminate shortcomings of that instrument and other actuarial approaches. The most important conclusion is that an instrument as good as the Static-99 generally cannot identify individuals who satisfy the legal requirements for sex offender civil commitment. By achieving effect sizes comparable to the Static-99 with more sensitive treatment of age and fewer variables, this Article provides further support for the view that the Static-99 does not properly account for age. By showing significant state effects using a model comparable to the Static-99, this Article underscores the importance of tailoring predictions to the particular jurisdiction. And, finally, the Article demonstrates in two different ways the large impact of prediction error on how many individuals will qualify for commitment.

¹⁴⁸ Woodworth & Kadane, *supra* note 85, at 239 (reporting somewhat better numbers in direct test of the MnSOST-R).

¹⁴⁹ Recall that this is equivalent to applying the preponderance of the evidence standard.

¹⁵⁰ As mentioned above, an alternative to this confidence interval approach to standards of proof is to set the cut-score in order to achieve a desired ratio of false positives ("FP") to false negatives ("FN"). Normally, these values are unavailable, Cohen, *supra* note 75, at 417, but here no one was civilly committed and we have actual data on recidivism. (Notably, this approach is independent of the commitment threshold and therefore probably not a good fit in this context.) The three standards of proof can be equated to FP:FN error ratios as follows: POE (50%), 1:1; CCE (75%), 1:3; and BRD (90%), 1:9. Applying these standards, 479, 346, and 149 individuals, respectively, qualified for commitment based on the logit predictions and observed recidivism. Thus, the unequal weighting of errors implied by heightened standards of proof can reduce commitments by up to 69%.

One criticism of the analysis above is that it is directed against a straw man: the Static-99 is not used in isolation. Rather, experts testify about the meaning of the score and offer opinions as to dangerousness that incorporate other factors. Existing data, however, suggest that adding clinical judgment to actuarial results does not improve predictive accuracy.¹⁵¹ Indeed, to the extent that there have been studies, they suggest that adjusting actuarial results actually decreases accuracy.¹⁵²

Another limitation is that recidivism information in these data is available only for the first three years after release. This generates a relatively low base rate. The observed sexual recidivism rate in the data is about 9.2%. In contrast, according to some researchers, “approximately 30% of sex offenders released from secure custody will have subsequent offenses recorded as sexual on police rap sheets.”¹⁵³ On the one hand, the low base rate makes the large effect sizes more impressive since prediction of low probability events is more difficult. This bolsters the present findings on age. But, on the other hand, the low base rate artificially reduces the number of individuals who qualified for commitment and perhaps exaggerates the impact of prediction error. One could respond by arguing that neither effect is “artificial.” When spending limited resources, the imminence of harm is plainly relevant. To prevent one sexual reoffense by locking up an individual for more than three years is arguably not cost-benefit justified. That point is, of course, debatable, and it must be conceded that the short follow-up period covered by the data is a limitation.

B. Implications

The way we select sex offenders for civil commitment is inadequate: essentially no one meets the legal standards. To combat this inadequacy, the practice of sex offender commitment should be curtailed or eliminated, the selection criteria lowered, or the selection methodology improved. The first two options involve policy judgments outside the scope of this Article.

¹⁵¹ See Terence W. Campbell & Gregory DeClue, *Flying Blind with Naked Factors: Problems and Pitfalls in Adjusted-Actuarial Sex-Offender Risk Assessment*, 2 OPEN ACCESS J. FORENSIC PSYCHOL. 75, 75 (2010), available at <http://www.forensicpsychologyunbound.ws/> (“Based on available data, at its best, [Adjusted Actuarial Assessment] neither increases nor decreases the accuracy of actuarial classification.”); Hamilton, *supra* note 16, at 44 (“[T]here is no empirical evidence that modifying actuarial scores improves the accuracy of predictions.”); see also Harris & Rice, *supra* note 46, at 1640.

¹⁵² See Campbell & DeClue, *supra* note 151, at 75 (“At its worst, [Adjusted Actuarial Assessment] dilutes actuarial accuracy.”); R. Karl Hanson & Kelly E. Morton-Bourgon, *The Accuracy of Recidivism Risk Assessments for Sexual Offenders: A Meta-Analysis of 118 Prediction Studies*, 21 PSYCHOL. ASSESSMENT 1, 7 (2009) (“[T]he adjusted scores showed lower predictive accuracy than did the unadjusted actuarial scores.”).

¹⁵³ Harris & Rice, *supra* note 46, at 1642. But see Hanson & Bussière, *supra* note 20, at 351 (reporting 13.4% sexual recidivism).

This Article does, however, shed light on a way forward for improved methodology. It should be emphasized that there is no guarantee that adopting one or even all of the suggestions below will solve the bottom-line problem.

Specifically, due to the low base rate of recidivism and substantial prediction error, an instrument as good as the Static-99 identified *not one* individual who qualified for commitment at the 50% or 75% threshold.¹⁵⁴ As noted above, half of the jurisdictions with sex offender commitment apply thresholds at or above 50%.¹⁵⁵ No evidence shows that adding other evidence to actuarial results improves prediction accuracy.¹⁵⁶ The obvious implication is that no one in these jurisdictions deserved to be civilly committed as a sex offender.

There are several responses to this finding. First, it depends crucially on the short follow-up period and resulting low base rate.¹⁵⁷ Arrest reports listing sexual offenses may understate recidivism for other reasons as well—e.g., failure of victims to report and failure of police to list the more difficult to prove sexual component of offenses like assault.¹⁵⁸ Second, the other half of jurisdictions have lower or unspecified commitment thresholds. This Article finds that standard could have been met with requisite certainty for a significant fraction of the sample. The problem, therefore, could be described as setting the threshold too high, rather than as failing to meet an impossible standard.

Still, this Article represents one of the first empirical tests of whether an instrument like the Static-99 can identify qualified individuals. The instrument failed to do so in half of jurisdictions. The base rate may be wrong or those jurisdictions may have the wrong standard, but it would seem that the burden going forward should be on the developers of ARAIs like the Static-99 to show that the instruments as revised can identify individuals who meet the commitment threshold at the required standard of proof. If no such showing is forthcoming, the entire enterprise of sex offender commitment is justifiably in doubt.¹⁵⁹

The Static-99 can be improved, at least in accounting for age, adjusting for jurisdiction-specific norms, and reporting error. The developers of the

¹⁵⁴ This finding stands in contrast to that of Janus and Meehl, who concluded as a theoretical matter that those standards could be met. Janus & Meehl, *supra* note 15, at 33.

¹⁵⁵ See *supra* Table 3.

¹⁵⁶ See *supra* notes 151–52 and accompanying text.

¹⁵⁷ Observed base rates vary as widely as 7.5% to 66.7%, with two large meta-analyses finding rates around 13–14% for five-year recidivism. Prentky et al., *supra* note 8, at 373–74.

¹⁵⁸ Harris & Rice, *supra* note 46, at 1643–44.

¹⁵⁹ Alternatives like longer prison sentences, or supervision and treatment in the community, may be preferred. See Cantone, *supra* note 109, at 720–21 (advocating for longer prison sentences); Eric S. Janus, *Minnesota's Sex Offender Commitment Program: Would an Empirically-Based Prevention Policy Be More Effective?*, 29 WM. MITCHELL L. REV. 1083, 1132–33 (2003) (advocating for stronger community-based treatment solutions).

Static-99 have recognized some of these shortcomings and offered updated alternatives. Recall that a revised version of the Static-99 includes four age categories instead of two.¹⁶⁰ This is certainly a step in the right direction, but why not go all the way and include age as a continuous variable? By similar token, new norms are necessary—as this Article confirms¹⁶¹—but including dummy variables for each jurisdiction, updating data each year, and reestimating the logit model outlined above has the potential to seamlessly adjust predictions as observed behavioral changes.¹⁶² Because one can directly calculate individual errors using this approach, a successful challenge along the lines of *Rosado* would be much less likely. Reporting such errors provides a critical link between risk assessment instruments and commitment decisions. In short, a logistic regression-based approach holds more promise than traditional actuarial methods.¹⁶³

Although this Article focused on the Static-99 and sex offender civil commitment, the lessons apply more generally to other actuarial instruments used in other contexts. There are many such contexts. Take, for example, the following list of criminal applications:

From the use of the IRS Discriminant Index Function to predict potential tax evasion and identify which tax returns to audit, to the use of drug-courier and racial profiles to identify suspects to search at airports, on the highways, and on city streets, to the use of risk-assessment instruments to determine pretrial detention, length of criminal sentences, prison classification, and parole eligibility, prediction instruments increasingly determine individual outcomes in our policing, law enforcement, and punishment practices.¹⁶⁴

The IRS formula is apparently based on regression analysis.¹⁶⁵ In contrast, drug-courier profiles have never been empirically validated.¹⁶⁶

The most commonly used risk assessment instrument in this country is

¹⁶⁰ Helmus et al., *supra* note 17.

¹⁶¹ See *supra* Section III.C.2.

¹⁶² See Woodworth & Kadane, *supra* note 85, at 238, 239 (advocating for “a standardized data base” and presuming that “prediction models will be developed and updated via logistic regression”).

¹⁶³ But see Harris & Rice, *supra* note 46, at 1639 (noting that “regression weights are unstable on replication”).

¹⁶⁴ HARCOURT, *supra* note 16, at 2.

¹⁶⁵ Bernard E. Harcourt, *From the Ne'er-Do-Well to the Criminal History Category: The Refinement of the Actuarial Model in Criminal Law*, 66 LAW & CONTEMP. PROBS. 99, 147 n.204 (2003).

¹⁶⁶ Morgan Cloud, *Search and Seizure by the Numbers: The Drug Courier Profile and Judicial Review of Investigative Formulas*, 65 B.U. L. REV. 843, 845 (1985).

the Level of Services Inventory Revised (LSI-R).¹⁶⁷ The LSI-R, used for parole and other purposes, is more extensive but closely analogous in structure to the Static-99.¹⁶⁸ The LSI-R includes a dummy variable based on age at first arrest, but not age at release.¹⁶⁹ If the goal is predicting recidivism, failing to include both as continuous variables is a mistake.¹⁷⁰ Failing to weight the items using regression analysis is another defect shared by the LSI-R. And, finally, reporting logit predictions along with errors could lead to better decision-making than, as the LSI-R does, merely lumping individuals into “low,” “medium,” and “high” risk groups.¹⁷¹ The LSI-R demonstrates that the present examination of the Static-99 has potentially broad importance.

V. CONCLUSION

The Static-99 is the most thoroughly researched tool for predicting sexual recidivism.¹⁷² Almost no one before this Article, however, empirically assessed the most critical question: can the Static-99 predict recidivism well enough to meet the legal standard for sex offender commitment?¹⁷³ The answer is mixed and qualified, but largely negative. The limitations of this study preclude any strong conclusions, but the findings at least suggest that the goals and methods of sex offender civil commitment need to be reevaluated. In the meantime, this Article identifies several ways in which the Static-99 and like instruments are deficient and can, and should, be improved.

¹⁶⁷ TRACY W. PETERS & ROGER K. WARREN, NAT'L CTR. FOR STATE COURTS, GETTING SMARTER ABOUT SENTENCING: NCSC'S SENTENCING REFORM SURVEY 17 (2006), available at sentencing.nj.gov/downloads/pdf/articles/2006/Aug2006/document09.pdf.

¹⁶⁸ See JAMES AUSTIN ET AL., RELIABILITY AND VALIDITY STUDY OF THE LSI-R RISK ASSESSMENT INSTRUMENT i (2003), available at <http://www.portal.state.pa.us/portal/server.pt> (search “LSI-R RISK”).

¹⁶⁹ *Id.* at 14 tbl.7.

¹⁷⁰ Indeed, a validation study of the LSI-R found the variable “Arrested under age 16” to be significant in predicting recidivism. *Id.* at 18.

¹⁷¹ A recent examination of the LSI-R found recidivism predictive power (AUC = 0.66 and 0.73) comparable to that achieved by this Article's main logit model (AUC = 0.66; *supra* Table 5). Sarah M. Manchak et al., *Utility of the Revised Level of Service Inventory (LSI-R) in Predicting Recidivism After Long-Term Incarceration*, 32 LAW & HUM. BEHAV. 477, 482 (2008).

¹⁷² See, e.g., Hanson & Morton-Bourgon, *supra* note 152, at 17 tbl.A1 (listing sixty-three such studies).

¹⁷³ Again, Hart et al., *supra* note 9, and Janus & Meehl, *supra* note 15, come closest.

APPENDIX

Formally, the logit model is specified as follows:

$$\text{Equation 1. } P_i = \frac{1}{1 + e^{-\beta X_i}}$$

where P_i is the probability of an individual hit or miss, e is the base of natural logarithms, β is a matrix of coefficients, and X_i a matrix of individual-specific variable values.¹⁷⁴

Two post-estimation calculations are complex enough to require explanation. In Table 7, I estimate the number of individuals who met the legal standards for commitment—for example, whose estimated likelihood of recidivism was above 50% (“more likely than not”) with 75% confidence (CCE). This required calculating the lower CI for logit predictions, P_i . I substituted for βX_i in Equation 1 one side of the following formula:

$$\text{Equation 2. } LB_p = LP_i - (Z_p \times SE_{LP_i})$$

where LB_p is the lower-bound of the linear CI for a given proof standard (75% or 90%), LP_i is the linear prediction of the logit model for an individual, Z_p is the inverse cumulative standard normal distribution for either 75% or 90%, and SE_{LP_i} is the standard error of LP_i .¹⁷⁵

The Cohen’s d statistic is defined as $(M_1 - M_2)/S_w$, where M_1 is the mean of one group, M_2 is the mean of the comparison group, and S_w is the pooled-within standard deviation.¹⁷⁶ The complicated part of this equation is the last term:

$$\text{Equation 3. } S_w = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2}}$$

where n is group size and s is group standard deviation.¹⁷⁷ The CIs on

¹⁷⁴ ROBERT S. PINDYCK & DANIEL L. RUBINFELD, *ECONOMETRIC MODELS AND ECONOMIC FORECASTS* 258 (3d ed. 1991).

¹⁷⁵ See Mark Inlow, *Prediction Confidence Intervals After Logistic Regression*, STATA (Apr. 1999, rev. July 2007), <http://www.stata.com/support/faqs/stat/prep.html>. By deriving the confidence interval from the standard error, this methodology avoids one criticism leveled against Hart et al., *supra* note 9. See Harris et al., *supra* note 78, at 154 (“The appropriate statistic is standard error of measurement . . .”).

¹⁷⁶ Hanson & Morton-Bourgon, *supra* note 152, at 5. See generally *Effect Size*, WIKIPEDIA, http://en.wikipedia.org/wiki/Effect_size (last visited July 4, 2011).

¹⁷⁷ JACOB COHEN, *STATISTICAL POWER ANALYSIS FOR THE BEHAVIORAL SCIENCES* 67 equation 2.5.2 (2d ed. 1988); Will Thalheimer & Samantha Cook, “How To Calculate Effect Sizes from Published Research: A Simplified Methodology,” 4 equation 1a, *available at*

Cohen's d statistics were calculated with the METAN downloadable add-on to Stata.¹⁷⁸ All computations in this Article were performed using Stata/SE 11.1.

docs.docstoc.com/pdf/.../6af10ee0-3d03-46ac-bd77-6a17477830e7.pdf (last visited Sept. 14, 2011); see also *Effect Size*, *supra* note 176.

¹⁷⁸ See Ross Harris et al., IDEAS, *METAN: Stata Module for Fixed and Random Effects Meta-Analysis*, <http://ideas.repec.org/c/boc/bocode/s456798.html> (last visited Sept. 26, 2011).

