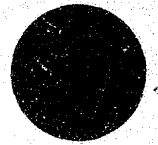# Identifying High-Rate Serious Criminals from Official Records

John E. Rolph, Jan M. Chaiken

The research described in this report was supported by the National Institute of Justice, U.S. Department of Justice under Grant No. 83-IJ-CX-0006.

R-3433-NIJ

# Identifying High-Rate Serious Criminals from Official Records

John E. Rolph, Jan M. Chaiken

April 1987

NCJRS

MAY 29 1987

ACQUISITIONS

# RAND

# PREFACE

This report documents one of a series of studies that are based on data collected in RAND's Second Inmate Survey, a project supported by the National Institute of Justice. The present work was funded under the Prediction, Classification, and Methodology program of research at the National Institute of Justice. The report should interest researchers and policymakers who care about the use of official record information and offender self-report information in distinguishing between offenders who commit serious crimes at high rates and those who do not.

Other reports concerned with the RAND Second Inmate Survey include:

1. Mark Peterson, Jan Chaiken, Patricia Ebener, and Paul Honig, *Survey of Prison and Jail Inmates: Background and Method*, N-1635-NIJ, November 1982. Describes the purposes of the survey, its design and administration, the data collected, and response patterns.

2. Kent Marquis with Patricia Ebener, *Quality of Prisoner Self-Reports: Arrest and Conviction Response Errors*, R-2637-DOJ, March 1981. Analyzes the reliability of the survey's self-reported arrest and conviction data, using both the retest method and a comparison with official records.

3. Jan Chaiken and Marcia Chaiken, with Joyce Peterson, *Varieties of Criminal Behavior: Summary and Policy Implications*, R-2814/1-NIJ, August 1982. Gives conclusions from analyses of the survey and official record data concerning the identification of serious criminal offenders and the implications of their behavioral characteristics for public policy.

4. Jan Chaiken and Marcia Chaiken, *Varieties of Criminal Behavior*, R-2814-NIJ, August 1982. Identifies ten subgroups of offenders and describes their behavioral characteristics, with special reference to the most serious offenders. Shows how, and the extent to which, serious offenders and high-crime-rate offenders can be identified from their characteristics and criminal records. Appendixes describe (a) an analysis of the internal consistency of survey responses and their correspondence with official record data, and (b) the construction of scaled predictor variables.

5. Peter W. Greenwood with Allan Abrahamse, *Selective Incapacitation*, R-2815-NIJ, August 1982.

6. Joan Petersilia, Paul Honig, with Charles Hubay, *The Prison Experience of Career Criminals*, R-2511-DOJ, May 1980. Describes the treatment need and program participation rates of prison inmates.

7. John E. Rolph, Jan M. Chaiken, and Robert L. Houchens, *Methods for Estimating Crime Rates of Individuals*, R-2730-NIJ, March 1981.

8. Stephen P. Klein and Michael N. Caggiano, *The Prevalence, Predictability and Policy Implications of Recidivism*, R-3413-BJS, August 1986.

John Rolph is on the research staff at The RAND Corporation. Jan Chaiken was on the research staff of The RAND Corporation during the early and middle phases of this study. He completed his work on the study as a consultant to The RAND Corporation. He is currently deputy manager of the law and justice research area at Abt Associates Inc.

# EXECUTIVE SUMMARY

Studies revealing the existence of a small group of criminal offenders who commit serious crimes at exceptionally high rates have inspired numerous efforts to identify these offenders and handle them differently from others. Achieving these goals could reduce crime rates or at least improve the efficiency of the criminal justice system. Yet a large gap separates the recognition from the identification. Criminal justice practitioners must answer many questions before they can confidently implement policies for dealing with high-rate offenders:

1.  *What is the relationship between high-rate serious offenders as found in self-report studies, and offenders who have lengthy official criminal histories or high rates of arrest?*

The 1970s notion of a "career criminal" was a person who commits crimes over an extended period of time and consequently has a lengthy official record of criminal activity. Of course, there may be substantial overlap between these offenders and high-rate serious offenders: Some people who have a record of criminal activity over an extended period are still committing crimes frequently. But many high-rate serious offenders have no official record, or minor records. And substantial numbers of offenders are arrested frequently without committing crimes at high rates. So the duration of a person's criminal record and the number of arrests he incurred over his career are not sure signs that he is a high-rate serious offender.

2.  *What is the relationship between recidivism and high-rate serious offending?*

By every indication, high-rate serious offenders are much more likely than the average offender to recidivate after incarceration (but the research evidence for this link is not very solid). The converse proposition, that recidivists are likely to be high-rate serious offenders, cannot be correct. The high-rate offenders constitute only a small fraction of any offender population, but typically more than half of an incarcerated population will recidivate. Many prediction instruments accurately identify recidivists, but they may or may not identify high-rate serious offenders. Recidivists could be predominantly offenders who commit less serious crimes or commit crimes at lower rates.

3.  *What levels of predictive accuracy are achieved by formulas or scales that have been suggested for use in identifying high-rate serious offenders?*

The best instruments developed for predicting either recidivism or "success" on bail release have impressive demonstrated levels of predictive validity. The same cannot be said of instruments proposed for identifying high-rate serious offenders. Until recently none of them has been tested as *predictive* instruments, in the sense of providing information about *future* commission of serious crimes at high rates. Studies of high-rate offenders have been retrospective: An offender's own reports about his past crime-committing behavior have been compared with other information about his past (official records, personal characteristics, drug use, employment, and the like). Recent exceptions to this practice are Greenwood and Turner (1987) and Klein and Caggiano (1986).

No previous studies have defined the concept of "high-rate *serious* offender" and developed discriminant rules that purport to distinguish between offenders who belong to this

category and offenders who do not. Some studies have attempted to estimate the person's crime commission rate for a particular type of crime—e.g., robbery—as a function of the person's other characteristics. Such equations may possibly have value for distinguishing high-rate offenders from others; but they may also be simply distinguishing very low-rate offenders from moderately low-rate offenders, or low-rate offenders from average offenders. Moreover, estimating an individual's commission rate for a single crime type does not capture the great variety of criminal behavior engaged in by most high-rate serious offenders.

Other studies have defined categories of serious offenders according to the combinations of crimes they commit. For example, an offender who simultaneously (over a one-to-two year period) deals drugs and commits both robbery and assault has been dubbed a "violent predator," and methods for discriminating between a person being a violent predator or not have been proposed (Chaiken and Chaiken, 1982). Even though the so-defined violent predators are much more likely than other offenders to be high-rate serious offenders, an equation or rule that identifies violent predators strictly speaking misses the mark of identifying high-rate serious offenders.

4. *Are typical official records sufficiently informative to be used in identifying high-rate serious offenders?*

Because studies of high-rate serious offenders are based on their self-reports, many of the data available in those studies, even if apparently describing information that should be reflected in official records, are often also obtained from self-reports. Such items as the person's lifetime count of felony arrests, whether he has ever been imprisoned, whether he was convicted of a crime before age 16, and whether he is addicted to drugs are typically also self-reports.

5. *Have the research reports that reveal prodigious levels of criminal activity by high-rate offenders actually exaggerated their true behavior?*

Criticisms have been raised about the mathematical methods used for estimating the annual crime commission rates of offenders from their self-reports (Visher, 1986). Moreover, the manner in which the resulting numbers have been displayed and summarized may tend to mislead nontechnical readers about the amount of crime committed by high-rate offenders. For example, it is mathematically correct to say that a person who committed eight robberies in two months (and was then arrested) was committing robberies at an annual rate of 48 robberies per year; but this statement tends to leave the possibly incorrect impression that the offender would indeed have committed 48 robberies if left free for an entire year.

Although action-oriented criminal justice practitioners may believe (or wish) that these questions were already settled or have obvious answers, techniques for identifying high-rate offenders are in their infancy and require further development. How much crime can be reduced by selectively focusing criminal justice system resources on high-rate offenders cannot be reliably estimated without having accurate estimates of crime commission rates for high-rate offenders. Selective programs cannot work as expected if they target the wrong offenders. For example, superficially similar people may be either low-rate or high-rate offenders.

## PURPOSE OF THIS STUDY

We set out to develop better ways of using official criminal records to conclude that certain offenders are committing crimes at high rates. Perhaps rules based on total adult convictions for specified crimes are too imprecise, or they look too far back into a person's past. Granted, some offenders who commit numerous crimes have unremarkable records (and cannot be identified from their records). But a record that shows extensive recent activity—for example, two burglary convictions, an assault conviction, and three arrests for robbery in the past year—might well indicate high-rate serious behavior. Two aspects of this example reflect the emphasis of this research: Information about several crime types is brought together, and arrests unconfirmed by convictions are taken into account if they are recent. Our study helps clarify *how, and in what circumstances, official records of arrests and convictions can help identify the serious offender.*

To address this and related issues, our study focuses on the offender who has a high annualized crime commission rate (number of crimes per year of free time). We pay special attention to offenders who commit *serious* crimes at high rates and to those who commit *multiple types of crimes* at high rates. If these offenders' criminal activities can be interrupted or terminated, the level of crime in society could potentially drop significantly.

In this study we have attempted to develop operational definitions of what constitutes a high-rate serious criminal. Then using statistical methods appropriate to discriminating between two possible outcomes, we assess how reliably official record information allows offenders to be labeled high-rate serious.

## METHODS

In this study we exploit the RAND Second Inmate Survey, which contains both self-reports of crime commissions and official records of arrests and convictions for several individuals. Although the data have already been subjected to an extensive amount of analysis for validity, the characteristics of serious criminal offenders, and the implications of various sentencing policies for selective incapacitation (Chaiken and Chaiken, 1982; Chaiken and Chaiken, 1985; Chaiken, Chaiken, and Peterson, 1982; Greenwood and Abrahamse, 1982; Greenwood and Turner, 1987; Marquis and Ebener, 1981; Visher, 1986), a few important issues remain about high-rate offenders on which we have tantalizing but imprecise findings.

Working with the same self-report data concerning incarcerated offenders that we use, Chaiken and Chaiken (1982) showed that straightforward approaches to identifying serious offenders from their official records lead to unsatisfactory—one might even say discouraging—results. They reported findings from two methods:

1. Categorizing offenders into so-called "complexes" according to the types of crimes they *commit* (without regard to their annualized rate of commission) and attempting to discriminate among complexes using official record information about present and past conviction offenses and drug use.

2. Carrying out multiple regression analyses, with the logarithm of self-reported annualized crime rate during the measurement period as the outcome variable, and official record data concerning current conviction offenses, prior convictions, and arrests during the measurement period as explanatory variables. This method attempts to distinguish among offenders according to the rates at which they commit crimes, disregarding the question of what particular crime or combination of crimes they commit.

With the first approach, the discrimination between the most serious complex of offenders (those who, by definition, committed robbery, assault, and drug deals in the measurement period—most of them having also committed various other crimes) and others was statistically significant, but the errors in discrimination were so extensive as to make the result pragmatically uninteresting. The multiple-regression analyses, in which Chaiken and Chaiken (1982) attempted to develop prediction equations for the rate at which offenders commit crimes, were even less promising. For survey respondents who had been convicted of robbery, for example, a regression equation using predictor variables based on the official record items listed above accounted for only 17 percent of the variance in logarithm of the annualized robbery rate.[1] Moreover, no convictions or arrests for crimes other than robbery entered into the regression[2] for the (logarithm of) robbery commission rate, although high-rate robbers were shown to be characterized by involvement in a wide variety of criminal activities.

These negative findings serve as a starting point for our research. The method of "complexes" does not use annualized commission rates (other than zero vs. not zero), and the regression methods are based on assuming smooth relationships across the entire range of values of annualized rates. Our approach considers all offenders whose reported commission rate for a particular crime type is above a certain threshold to be indistinguishable in regard to that type of crime. Specifically, we investigate various definitions of whether an offender is "high-rate serious." We develop a different method of estimating an individual's annual crime commission rates[3] than Chaiken and Chaiken (1982) used, and we establish threshold values for labeling an offender "high rate." Our alternative definitions of high-rate serious specify which crime types, committed alone at high rates or in combination with certain other crimes at high rates, constitute the kind of behavior to be analyzed. These outcome variables for our analysis combine both concepts of serious offending used in previous research: committing specified combinations of crimes, and committing one or more types of crimes at high rates.

Our first step in developing an operational definition of high-rate serious offender is to estimate individual offenders' self-reported crime commission rates for the various crimes. Chaiken and Chaiken (1982) and other workers following them used the "annualized crime commission rate" defined in the obvious way: the number of reported crimes of that type divided by the amount of unincarcerated time during the measurement period. By this definition, some respondents with extremely high estimated rates have very short periods of unincarcerated time during their measurement periods. One might argue that an offender who displayed a prodigious rate of criminal activity for a short period of time (say one to four months) would not (or could not) sustain this rate for an entire year. If so, the calculated "annualized crime commission rate" overrepresents the number of crimes he would commit if left unincarcerated for a year. There is evidence of "spurting" behavior—committing crimes at high rates for short periods of time. Therefore we developed an adjusted estimate of each offender's annual crime commission rate that takes into account the variation throughout the year of an individual offender's crime commission rate for a particular crime. The result of this adjustment is to reduce the more extreme annualized rates considerably.

The second step in developing operational definitions of high-rate serious offenders was exploring sensible ways of specifying which crime types, committed alone or in combination,

---

[1]When self-report items similar to the official record items were used in the regression, substantially larger percentages of the variance were explained.

[2]Using conventional statistical tests of significance of the estimated coefficients.

[3]Rolph, Chaiken, and Houchens (1981) analyzed the RAND Second Inmate Survey data and found some indication that the unadjusted annualized rates may be extreme for high offenders. We use an adjusted rate here.

should be considered as constituting *serious* behavior. Given a definition of serious, all offenders whose adjusted annual rate estimates are above a certain threshold for those crimes are deemed *high-rate serious*. An individual is defined to be a "high-rate serious" offender in each of the six interpretations as follows:

1. If he has a high robbery commission rate.
2. If he has a high robbery rate or a high rate of robbery of persons or a high rate of robbery of businesses.
3. If interpretation 2 applies or he has a high rate of committing assault.
4. If *any* of his crime rates are high[4] and he is a violent predator (i.e., commits robbery and assault and deals drugs).
5. If any of his crime rates are high and he commits robbery and assault.
6. If any of his crime rates are high and he deals drugs and commits either robbery or assault.

A very useful early finding from the analysis of these six interpretations was that it made little difference whether the adjusted or unadjusted crime rates were used in the definitions. This means that to a large extent offenders in the data who were classified as high-rate serious according to, say, interpretation 4, using the 80th percentile of the unadjusted rate as the cut-off for high rates, were also classified as high-rate serious using the 80th percentile of the adjusted rate. Correspondingly, few offenders had different classifications using the unadjusted 80th percentile and the adjusted 80th percentile.

We develop logistic regression equations for estimating the probability that an offender is high-rate serious (using a particular definition) as a function of his official record characteristics (personal characteristics, background factors, arrests, convictions, etc.). The fitted logistic equations are used to define a discriminant rule of the form: Classify individual n as high rate if his estimated probability $p(x_n)$ is above a certain threshold value (where $x_n$ is the vector of person n's characteristics). For comparison purposes, discriminant rules for identifying high-rate offenders are also developed using *self-report* versions of official record characteristics as the explanatory variables. A more complete set of explanatory variables based on all available self-report information is also used to estimate discriminant rules. We evaluate these various discriminant rules by computing their error rates in correctly labeling which offenders are in fact high-rate serious as defined by the outcome variable used in the analysis.

The population base for fitting logistic regression equations consists of all prisoner respondents to RAND's Second Inmate Survey who satisfied two conditions:

- Their official record data (inmate folders) had been located and coded by the original researchers.
- They reported committing one or more of the ten crimes included in the survey instrument.[5]
- Respondents to the Second Inmate Survey who were surveyed in jails rather than prisons were excluded from the present study because no official record information was collected for them in the original study. Respondents who said they committed none of the crimes in the survey booklet were excluded because they may have been high-

---

[4]That is, for any of the types of crimes studied in the RAND Second Inmate Survey.

[5]The tabulations of adjusted crime commission rates, presented in App. A, include respondents who denied committing all of the crimes included in the survey booklet, but they are limited to prisoner respondents for whom official record data were originally collected.

rate or low-rate in regard to the crimes they do commit, but the data would not permit making the necessary distinctions. (The survey booklet was not comprehensive in its coverage of different crime types, omitting such serious crimes as homicide, kidnap, and forcible rape.)

## RESULTS

We were somewhat surprised and disappointed that none of these definitions appeared to capture a natural division between the "really bad guys" and the other offenders. Indeed, all of our definitions of "high-rate serious offender" had similar properties with respect to which explanatory variables discriminated between them. After we excluded two definitions that yielded too few high-rate serious offenders for reliable statistical analyses, our preferred definition was a somewhat arbitrary choice.

We had an unsuccessful search for discriminant rules based on an offender's official record information that would reliably label high-rate serious offenders correctly. Our best discriminant rules based on official record information did about 20 percent better than a "chance rule" in correctly labeling such offenders using our preferred definition. And the explanatory variables capturing the official record arrest and conviction information available in the RAND survey has a modest but statistically significant relationship with being a high-rate serious offender in each of the three states. This experience is consistent with the results of Chaiken and Chaiken (1982), who did not use statistical methods specifically tailored to the discrimination problem.

When we use the self-report versions of the official record explanatory variables in an attempt to overcome the limitations of possibly low quality available official record data, the situation improves slightly. This discriminant rule correctly labels high-rate serious offenders 25 percent better than a chance rule. This finding is consistent with those of Chaiken and Chaiken (1982), who detected a slight improvement in estimating robbery rates when using self-report versions of official record information over the official record information.

Finally, as a benchmark we develop our "best" discriminant rule based on all self-report data. The improvement in correctly identifying high-rate serious offenders is substantial— almost 40 percent better than a chance rule. The improvement is due in large part to variables capturing aspects of the offender's juvenile period (crime, commitment to state facility, heroin abuse, and high school graduation) and the offender's social circumstances (employment, substance abuse) during the measurement period. Including explanatory variables on the offender's race in addition to the above variables did not improve the correct identification rate of the discriminant rule. Although the precise performance measure values of our discriminant rule vary with the definition of high-rate serious offender, the patterns are as described above.

## POLICY CONCLUSIONS

This research was motivated in part by the debate surrounding selective incapacitation (Blackmore and Welsh, 1984; Chaiken and Chaiken, 1985; Cohen, 1983; Fischer, 1984a; Forst, 1983; Greenwood and Abrahamse, 1982; Greenwood and Turner, 1987; Spelman, 1986; von Hirsch, 1984, 1985; von Hirsch and Gottfredson, 1984). Our research does not address the question of how predictive past behavior is of future commissions of crime. Prospective prediction is required for a direct test of any selective incapacitation sentencing, probation, or parole

policy. Greenwood and Turner (1987) and Klein and Caggiano (1986) have carried out studies of this nature. Howevei, if one can assume that trends in offenders' crime commission rates change only slowly over time, examining the relationship between concurrent arrests or convictions and self-reported crime commissions is relevant to prediction of future offenses.

Some workers have achieved respectable power with discriminant rules aimed at predicting recidivism (e.g., Fischer, 1984a); we therefore expected better success than we actually achieved in distinguishing high-rate serious offenders using official record data. We found no evidence that the relationship in the RAND Second Inmate Survey between available official record variables, or indeed any set of explanatory variables, and being a high-rate serious offender is strong enough to be of practical use for many criminal justice policy purposes. However the substantial improvement that our "full self-report" model rule made in discriminatory power gives some promise that using discriminant rules with carefully recorded juvenile record information may improve identification of high-rate serious offenders.

# ACKNOWLEDGMENTS

# CONTENTS

# FIGURES

# TABLES

# I. INTRODUCTION

What are good ways to design programs aimed specifically at career criminals? People will answer this question according to their implicit definition of the term "career criminal," the ways in which they think career criminals can be identified, their objectives in dealing with criminals, and the type of program they have in mind—employment, drug treatment, intensive prosecution, longer sentences for those convicted, etc. (Blumstein et al., 1986; Chelimsky and Dahmann, 1981). Definitions color one's impressions of the types of criminals who fit the label "career criminal" and thus suggest what should be done about handling them.

Today's notion that a career criminal is an offender who commits a disproportionately large number of crimes reflects a substantial change from earlier usage of the term "career," which connoted someone who commits crime over an extended period of time (Petersilia and Lavin, 1978). Of course, there is substantial overlap: Some people who have a record of criminal activity over an extended period are still committing crimes frequently. But the shift in definitions has come about because many offenders with long criminal records are no longer committing crimes at high rates, while others with minor or nonexistent records are highly active offenders (Chaiken and Chaiken, 1982; Chaiken and Chaiken, 1985; Greenwood and Turner, 1987; Klein and Caggiano, 1986; Petersilia, 1980; Peterson and Braiker, 1981; Rhodes et al., 1982).

Many criminal justice practitioners, motivated by their own experience and research knowledge, are moving away from policies that define career criminals in terms of lengthy official criminal histories (Forst, 1982). But the irony of this development is that older definitions were easy to formulate and implement, they had an appearance of fairness, and in some cases they worked. By contrast, attempting to focus on the offender who commits a disproportionately large number of crimes is fraught with pitfalls because it is difficult, if not impossible, to find an accurate method of identifying such offenders.

## SOME RESEARCH QUESTIONS

Studies revealing the existence of a small group of criminal offenders who commit serious crimes at exceptionally high rates have inspired numerous efforts to identify these offenders and handle them differently from others (Chaiken and Chaiken, 1984; Elliott, Ageton, and Huizinga, 1980; Peterson and Braiker, 1981; Tracy, Wolfgang, and Figlio, 1985; Wolfgang, Figlio, and Sellin, 1972). Achieving these goals could reduce crime rates or at least improve the efficiency of the criminal justice system. Yet a large gap separates the recognition that high-rate serious offenders exist from the ability to identify them. Many questions must be answered before criminal justice practitioners can confidently implement policies for dealing with high-rate offenders:

1.  *What is the relationship between high-rate serious offenders, as found in self-report studies, and offenders who have lengthy official criminal histories or high rates of arrest?*

Although some people who have a lengthy record of criminal activity are high-rate offenders, others are merely "losers" who commit few crimes but are arrested often (Chaiken and Chaiken, 1985). In addition, some high-rate serious offenders have no official record, or

1

only minor records. So the duration of a person's criminal record or recent high rates of arrest are not sure signs that he is a high-rate serious offender.

### 2. What is the relationship between recidivism and high-rate serious offending?

By every indication, high-rate serious offenders are much more likely than the average offender to return to criminal activity after being incarcerated (but the research evidence for this link is not very solid). The converse proposition, that recidivists are likely to be high-rate serious offenders, is not correct. Typically more than half of an incarcerated population will recidivate, so the number of recidivists is much larger than the number of high-rate serious offenders. Many prediction instruments accurately identify recidivists (Fischer, 1984b; Gottfredson and Gottfredson, 1986; Hoffman, 1983), but they may or may not identify high-rate serious offenders. Recidivists could be predominantly offenders who commit less serious crimes or commit crimes at lower rates.

### 3. Are typical official records sufficiently informative to be used in identifying high-rate serious offenders?

Because studies of high-rate serious offenders are based on their self-reports, many of the data available in those studies, even if apparently describing information that should be reflected in official records, are often also obtained from self-reports. Thus such items as the person's lifetime count of felony arrests, whether he has ever been imprisoned, whether he was convicted of a crime before age 16, and whether he is addicted to drugs are also typically self-reports in those studies. The predictive value of information actually collected from official records is very low compared with the strength of self-report versions of data that might also be in official records (Chaiken and Chaiken, 1982).

### 4. What levels of predictive accuracy are achieved by formulas or scales that have been suggested for use in identifying high-rate serious offenders?

The best instruments developed for predicting either recidivism or "success" on bail release have impressive, demonstrated levels of predictive validity. The same cannot be said of instruments proposed for identifying high-rate serious offenders. Until recently none of them has been tested as *predictive* instruments, in the sense of providing information about *future* commission of serious crimes at high rates. That is, all studies of high-rate offenders have been retrospective: An offender's own reports about his past crime-committing behavior have been compared with other information about his past (official records, personal characteristics, drug use, employment, and the like). Recent exceptions to this practice are Greenwood and Turner (1987) and Klein and Caggiano (1986).

Before the start of the present study, no researchers had defined the concept of "high-rate *serious* offender" and developed discriminant rules that purport to distinguish between offenders who belong to this category and offenders who do not.[1] Some studies have attempted to estimate the person's crime commission rate for a particular type of crime—e.g., robbery—as

---

[1] The work of Chaiken and Chaiken (1985) carried out in parallel with the present study uses some of the same concepts.

a function of the person's other characteristics.[2] Such equations may possibly have value for distinguishing high-rate offenders from others; but, because they focus on the full range of variation in crime commission rates, they may also be simply distinguishing very low-rate offenders from moderately low-rate offenders, or low-rate offenders from average offenders. Moreover, estimating an individual's commission rate for a single crime type does not capture the great variety of criminal behavior most high-rate serious offenders engage in.

Other studies have defined categories of serious offenders according to the combinations of crimes they commit. For example, an offender who simultaneously (over a one-to-two year period) deals drugs and commits both robbery and assault has been dubbed a "violent predator," and methods for discriminating between a person's being a violent predator or not have been proposed (Chaiken and Chaiken, 1982). Even though the so-defined violent predators are much more likely than other offenders to be high-rate serious offenders, an equation or rule that identifies violent predators does not necessarily identify high-rate serious offenders.

5. *Have the research reports that reveal prodigious levels of criminal activity by high-rate offenders actually exaggerated their true behavior?*

Criticisms have been raised about the mathematical methods used for estimating the annual crime commission rates of offenders from their self-reports (Visher, 1986). Moreover, the manner in which the resulting numbers have been displayed and summarized may mislead nontechnical readers about the amount of crime committed by high-rate offenders. For example, it is mathematically correct to say that a person who committed eight robberies in two months (and was then arrested) was committing robberies at an annual rate of 48 per year; but this statement tends to leave the possibly incorrect impression that the offender would indeed have committed 48 robberies if left free for an entire year.

Although action-oriented criminal justice practitioners may believe (or wish) that these questions were already settled or have obvious answers, techniques for identifying high-rate offenders are actually in their infancy and require further development. How much crime can be reduced by selectively focusing criminal justice system resources on high-rate offenders cannot be reliably estimated without having accurate estimates of crime commission rates for high-rate offenders. Selective programs cannot work as expected if they target the wrong offenders. For example, superficially similar people may be either low-rate or high-rate offenders.

## PURPOSE OF THIS STUDY

We set out to develop better ways of using official criminal records to conclude that certain offenders are committing crimes at high rates. Perhaps rules based on total adult convictions for specified crimes are too imprecise, or they look too far back into a person's past. Granted, some offenders who commit numerous crimes have unremarkable records (and cannot be identified from them). But a record that shows extensive recent activity—for example, two prior burglary convictions, an assault conviction, and three arrests for robbery in the past year—might well indicate high-rate serious behavior. Two aspects of this example reflect the emphasis of this research: Information about several crime types is brought together, and

---

[2]The studies may estimate the crime commission rate itself, or a mathematical transformation such as the logarithm.

arrests unconfirmed by convictions are taken into account if they are recent. Our study helps clarify *how, and in what circumstances, official records of arrests and convictions can help identify the serious offender.*

To address this and related issues, our study focuses on the offender who has a high annualized crime commission rate (number of reported crimes divided by length of reporting period). We pay special attention to offenders who commit *serious* crimes at high rates and to those who commit *multiple types of crimes* at high rates. If these offenders' criminal activity can be interrupted or terminated, the level of crime in society could potentially drop significantly.

Studies of self-reported criminal activity show that, for any given type of crime, a small but significant fraction of the study population of criminals reports committing the crime at annualized rates as much as 50 times higher than the median rate of all those who commit the crime (Chaiken and Chaiken, 1982; Elliott, Ageton, and Huizinga, 1980; Peterson and Braiker, 1981; Reiss, 1973; Visher, 1986). In other words, the self-reports of some offenders suggest they are dramatically more active than other offenders. Aside from issues of the believability of such self-reports, which we address below, such individuals are prime candidates for the label "high-rate serious offender."

Similarly, official records of police contacts with offenders show that a small fraction of the offenders have total lifetime numbers of contacts, or annualized contact rates during a specified period of years, that are substantially larger than others (Blumstein and Cohen, 1980; Shannon, 1978; Tracy, Wolfgang, and Figlio, 1985; Wolfgang, Figlio, and Sellin, 1972). These are the offenders who are often referred to as having "high rates" or being "responsible for a disproportionately large part of the crime problem." But with current knowledge, the extent to which these offenders coincide with offenders who actually have high crime commission rates is unknown. We are not asking whether the two groups coincide exactly; it is obviously possible in principle for an offender to have an extensive criminal record and yet commit crimes at a relatively low rate, or commit crimes at a high rate but have a minor or no record. We are asking instead a more fundamental question: On the whole are the offenders with extensive records high-rate offenders?

This study develops operational definitions of what constitutes a high-rate serious criminal. Then, using statistical methods appropriate to discriminating between two possible outcomes, it assesses how reliably official record information allows offenders to be labeled "high-rate serious."

## METHODS

In this study we use the RAND Second Inmate Survey data, which include both self-reports of crime commissions and official records of arrests and convictions for a number of individuals. See Ebener (1983) and Honig (1983) for a description of the data files and codebook. Although the data have already been analyzed extensively for validity, the characteristics of serious criminal offenders, and the implications of various possible sentencing policies for selective incapacitation, a few important issues remain about high-rate offenders on which we have tantalizing but imprecise findings (Chaiken and Chaiken, 1982; Chaiken and Chaiken, 1985; Chaiken, Chaiken, and Peterson, 1982; Greenwood and Abrahamse, 1982; Greenwood and Turner, 1987; Marquis and Ebener, 1981; Spelman, 1986; Visher, 1986).

Chaiken and Chaiken (1982), working with the same self-report data concerning incarcerated offenders that we use, showed that straightforward approaches to identifying serious

offenders from their official records lead to unsatisfactory—one might even say discouraging—results. They reported findings from two methods:

1. Categorizing offenders into so-called "complexes" according to the types of crimes they *commit* (without regard to their annualized rate of commission) and attempting to discriminate among complexes using official record information about present and past conviction offenses and drug use.
2. Carrying out multiple regression analyses, with the logarithm of self-reported annualized crime rate during the measurement period as the outcome variable, and official record data concerning current conviction offenses, prior convictions, and arrests during the measurement period as explanatory variables. This method attempts to distinguish among offenders according to the rates at which they commit crimes, disregarding the question of what particular crime or combination of crimes they commit.

Using the first approach, the discrimination between the most serious complex of offenders (those who, by definition, committed robbery, assault, and drug deals in the measurement period—most of them having also committed various other crimes) and others was statistically significant, but the errors in discrimination were so extensive as to make the result pragmatically uninteresting.

The multiple-regression analyses, in which Chaiken and Chaiken (1982) attempted to develop prediction equations for the rate at which offenders commit crimes, were even less promising. For survey respondents who had been convicted of robbery, for example, a regression equation using predictor variables based on the official record items listed above accounted for only 17 percent of the variance in logarithm of the annualized robbery rate.[3] Moreover, no convictions or arrests for crimes other than robbery entered into the regression[4] for the (logarithm of) robbery commission rate, even though high-rate robbers were shown to be predominantly characterized by involvement in a wide variety of criminal activities.

These negative findings serve as a starting point for our research. The method of "complexes" does not use annualized commission rates (other than zero vs. not zero), and the regression methods are based on assuming smooth relationships across the entire range of values of annualized rates. Our approach considers all offenders whose reported commission rate for a particular crime type is above a certain threshold to be indistinguishable in regard to that type of crime. Specifically, we investigate a variety of definitions of whether an offender is "high-rate serious." We develop a different method of estimating an individual's annual crime commission rates[5] than Chaiken and Chaiken (1982) used, and we establish threshold values for labeling an offender "high rate." Our alternative definitions of high-rate serious specify which crime types, committed alone at high rates or in combination with certain other crimes at high rates, constitute the kind of behavior to be analyzed. These outcome variables for our analysis combine both concepts of serious offending used in previous research: committing specified combinations of crimes, and committing one or more types of crimes at high rates.

We develop logistic regression equations for estimating the probability that an offender is high-rate serious (using a particular definition) as a function of information on his official

---

[3]When self-report items similar to the official record items were used in the regression, substantially larger percentages of the variance were explained.

[4]Using conventional statistical tests of significance of the estimated coefficients.

[5]Rolph, Chaiken, and Houchens (1981) analyzed the RAND Second Inmate Survey data and found some indication that the unadjusted annualized rates may be extreme for high offenders. We use an adjusted rate here.

record (personal characteristics, background factors, arrests, convictions, etc.). The fitted logistic equations are used to define a discriminant rule of the form: Classify individual n as high rate if his estimated probability $p(x_n)$ is above a threshold value (where $x_n$ is the vector of n's characteristics). For comparison purposes, discriminant rules for identifying high-rate offenders are also developed using *self-report* official record characteristics as the explanatory variables. A more complete set of explanatory variables based on all available self-report information is also used to estimate discriminant rules. We evaluate these various discriminant rules by computing their error rates in correctly labeling which offenders are in fact high-rate serious as defined by the outcome variable used in the analysis.

The population base for fitting logistic regression equations consists of all prisoner respondents to RAND's Second Inmate Survey who satisfied two conditions:

- Their official record data (inmate folders) had been located and coded by the original researchers.
- They reported committing one or more of the ten crimes included in the survey instrument.[6]

Respondents to the Second Inmate Survey who were surveyed in jails rather than prisons were excluded from the present study because no official record information was collected for them in the original study. Respondents who said they committed none of the crimes in the survey booklet were excluded because they may have been high-rate or low-rate in regard to the crimes they do commit, but the data would not permit making the necessary distinctions. The survey booklet was not comprehensive in its coverage of different crime types, omitting such serious crimes as homicide, kidnap, and forcible rape.


## ORGANIZATION

Section II describes the process we followed in defining high-rate serious behavior. In particular it covers our method for adjusting the estimated annual crime rates to account for the differing measurement periods and the "spurting behavior" of offenders. It also describes how we arrived at some 12 different definitions of high-rate serious offenders and the choices among these we made for analysis.

Section III describes how we selected explanatory variables to use in the discriminant rules for identifying high-rate serious offenders. There are three sets of explanatory variables used in different logistic regression equations: official record variables, self-report versions of official record variables, and a complete best set of variables. The explanatory power of the three sets of variables is in the order listed above.

Section IV evaluates how well the discriminant rules based on the logistic regression equations were able to identify high-rate serious offenders. The discriminant rules based only on official record information performed fairly poorly in correctly labeling high-rate serious offenders. The more complete self-report information yielded a rule that performed much better than the official record information rules but still made many errors.

Section V contains our conclusions.

---

[6]The tabulations of adjusted crime commission rates, presented in Appendix A, include respondents who denied committing all of the crimes included in the survey booklet. They are, however, limited to prisoner respondents for whom official record data were originally collected.

# II. DEFINING HIGH-RATE SERIOUS BEHAVIOR

How can one know whether a particular person is or is not a high-rate serious offender? This question has at least two different interpretations. The first assumes information is available describing the person's *actual* criminal behavior during some period of time. Subsequent sections examine the same question assuming only indirect or presumably related information is available—such as the person's *official* criminal record or history of drug use.

Even when offenders have provided detailed self-reports concerning their criminal behavior, determining whether they are high-rate serious offenders is not entirely straightforward. Many judgments must be made about matters that are at best ambiguous and at worst subject to vigorous dispute. The basic steps in making this determination for a particular person are:

- Estimating his annual crime commission rate for each of several crime types from the data provided,
- Establishing a way of labeling levels of estimated crime rates as either "high" or "not high," and
- Deciding which of these crime types, committed alone or in combination with specified other crimes, should be considered as constituting "serious" behavior.

This section describes how we carried out each of these steps in turn to arrive at a definition of high-rate serious behavior.

## ESTIMATING CRIME COMMISSION RATES FROM SELF-REPORT DATA

When the present study began, methods for estimating annualized crime commission rates for individuals were clouded by considerable controversy. Even though RAND researchers had conducted a series of increasingly sophisticated surveys of incarcerated criminal offenders, together with associated methodological studies,[1] doubts remained as to the validity and interpretation of the crime commission rates produced in those studies (see Visher, 1986). This section discusses the major sources of debate and presents our findings on these issues, including suitable methods of revising the previously published crime rate estimates based on RAND Inmate Survey data. Appendix A contains new tabulations of estimated crime commission rates, based on the methods described here.

A large amount of judgment can enter into the derivation of estimated annual crime commission rates from self-reports of crime commissions. Ambiguities in the respondent's answers to questionnaire items provide one area for judgment to play a role. For example, the respondent may provide conflicting information, may specify a range (e.g., "3 to 5") where a number is requested, give a verbal instead of a numerical response ("most," "a lot"), or give incomplete rate data (e.g., checking a box on the questionnaire labeled "several times a week" but not filling in an answer to the associated question "how many times a week?").

A second area for judgment occurs if, as in the case of the RAND surveys, two different formats of questions about crime commissions are included on the survey instrument

---

[1]See Chaiken and Chaiken (1982), Marquis and Ebener (1981), Petersilia, Greenwood, and Lavin (1977), Peterson and Braiker (1981), Peterson et al. (1982), Rolph, Chaiken, and Houchens (1981).

(primarily to permit reliability checks). The estimates of crime commission rates resulting from analysis of data from the two sets of formats may differ systematically, in which case there is an issue whether to consider one estimate preferable to the other for some reason, or to combine the two estimates in some way.

Third, the survey's sample design may include known selection biases, which can be corrected by various statistical methods. For example, if a sample is chosen from among people who admit to committing robbery during the past year, then their answer to the question "how many robberies did you commit in the last year" will always be one or more (i.e., not zero). Although some of the respondents' rates of committing robbery may be less than one per year, this cannot be seen directly from the data because of the method of selecting the sample and asking the question. Judgment enters into the analyst's decision whether and how to adjust the resulting estimates to take into account the selection bias.

Finally, analysts may differ in their mental model of the meaning of a crime commission rate. For example, if they assume an individual's crime commission rate remains constant over periods of many months or years, they will derive a different estimate from the data than if they assume the rate fluctuates in some specified way. Even if they simply convert the data into annual crime rates (for example, 11 crimes committed in a period of four months is equivalent to an annual rate of $11 \times 12/4 = 33$ crimes per year), they give the impression that they believe the individual is able to sustain the same rate for an entire year. However, the purpose of converting to annual rates is often merely to obtain standard statistics that can be compared across studies.

The data collected in RAND's surveys of incarcerated male criminal offenders have provided an opportunity to explore these areas of judgment and to determine which of them seem to make a sizable difference when estimating crime commission rates. Both the formats of the survey instruments and the fact that the data have been analyzed by several independent groups of researchers have contributed to our understanding of the important issues related to estimating crime commission rates.

RAND's 1976 inmate survey was an anonymous written survey of California prisoners.[2] The 1978–79 survey covered three states (Michigan and Texas as well as California) and included inmates of county jails as well as state prisons.[3]

The questionnaire instrument used in the 1978–79 survey was a refinement and elaboration of the first instrument. But both included questions asking for information about the numbers of crimes respondents had committed in a period preceding their incarcerations.

The published distributions of estimated crime commission rates differed substantially between the two surveys. Estimated mean rates from the second survey ranged between 1.6 and 20 times as high as the estimated rates for the same crime in the first survey. (See Table B.1 in App. B.) The Second Inmate Survey instrument replicated selected 1976 survey questions about crime counts and thus permitted us to analyze the role of analytical decisions and instrumentation differences in creating the disparities between the results of the two surveys.

The primary explanations for the disparities were found to be differences in the *length of the measurement period* in the two surveys and the researchers' decisions on *editing and interpreting of data* from a small number of respondents whose answers were unclear but indicated possibly high crime commission rates. Other possible explanations that turned out to be less important are discussed in App. B.

---

[2]See Peterson and Braiker (1981).

[3]For a description of the survey, see Peterson et al. (1982). The California prisoners were surveyed in 1979.

### Comparing Two Questionnaires: The Effects of Duration of Measurement Period

The questions about crime commissions that were included in both the 1976 and the 1978–79 RAND surveys asked for categorized counts of crimes committed by the respondents during a specified period.[4] (See App. C for the 1976 version of the questions, and App. D for the 1978–79 replication.) For example, respondents could answer that during the measurement period they committed 1–2 burglaries, 3–5 burglaries, 6–10 burglaries, or more than ten burglaries.

To analyze the effects of the duration of the measurement period on estimates of crime commission rates, we compared the responses of California prisoners to five of these questions in the two surveys. The categorized counts of crimes committed, tabulated in Table 1, are not adjusted for the duration of the measurement period.[5]

The similarity between the 1976 and 1979 raw data is striking, especially when compared with the differences apparent in the estimated annual crime rates published from the two surveys (Table B.1 in App. B). For auto theft, the two distributions in Table 1 are nearly identical.

For cons, the distributions are the same in the ranges 0, 1–2, and 3–5; relatively more of the 1979 responses are in the "over 10" category (compared with those in 1976) and correspondingly less in the 6–10 category. For forgery, more 1979 respondents reported "zero" and also more reported "over 10," but overall the distributions are not very different.

Only for burglary and drug sales do we find a substantial upward shift from the 1976 respondents to the 1979 respondents. For these two crimes substantially fewer 1979 respondents answered "zero" and substantially more answered "over 10." Considering that the measurement period was three years long in the 1976 survey and one to two years long in the 1979 survey, an indisputably higher level of criminal activity was reported in 1979 in the areas of burglary and drug sales.

The results presented in Table 1 strongly challenge the notion that offenders' crime commission rates remain constant over periods as long as several years. The distributions of total counts of crimes reportedly committed by respondents to the 1976 survey are on the whole closely comparable to the corresponding distributions for California prisoner respondents in the 1978–79 survey. Thus a large part of the difference between estimated crime rates from the two surveys arises from the fact that the first survey's numbers were divided by "street times" as long as three years, while the "street time" for the second survey ranged from a minimum of one month, to a typical figure around 14 months, to a maximum of 24 months.

The results support the view that much of the respondent's criminal activity occurred near the time of his arrest. The total count of crimes reported increases very little as the measurement period increases from six months to two years to three years. However, the annualized crime rates differ by ratios of 1:4:6 in these three cases simply because of changes in the denominator. This model of criminal behavior is consistent with our earlier analysis of Second Inmate Survey data, showing that counts of crimes reportedly committed were not strongly related to the duration of the respondent's street time (Rolph, Chaiken, and

---

[4]Responses to these questions have not been used in any earlier published analyses of the crime commission rates of respondents to the 1978–79 RAND Inmate Survey. (They were used only for measuring the internal reliability of respondents' answers.) The previously published crime commission rates were based on answers to a different set of questions (an example of which is shown in App. E).

[5]The respondents are not weighted in the tabulations presented in Table 1. Similar tables published by Peterson and Braiker (1981) based on the 1976 survey may differ because they weighted the respondents to represent a simulated cohort of offenders sentenced to prison in a given year.

Table 1

COMPARISON OF NUMBER OF CRIMES REPORTED ON 1976 SURVEY
AND CORRESPONDING 1979 REPLICATION QUESTION,
CALIFORNIA PRISONERS, UNWEIGHTED[a]

| Crime Type | Count of Crimes Committed | Percent | | Cumulative Percent | |
|---|---|---|---|---|---|
| | | 1976 Survey | 1979 Survey | 1976 Survey | 1979 Survey |
| Burglary | 0 | 54.3 | 45.6 | 54.3 | 45.6 |
| | 1-2 | 15.8 | 15.1 | 70.1 | 60.8 |
| | 3-5 | 10.8 | 12.8 | 80.9 | 73.5 |
| | 6-10 | 7.9 | 7.0 | 88.8 | 80.5 |
| | over 10 | 11.2 | 19.5 | 100.0 | 100.0 |
| Auto theft | 0 | 72.7 | 75.0 | 72.7 | 75.0 |
| | 1-2 | 16.6 | 11.3 | 89.3 | 86.3 |
| | 3-5 | 5.1 | 5.5 | 94.5 | 91.9 |
| | 6-10 | 3.0 | 4.7 | 97.5 | 96.5 |
| | over 10 | 2.5 | 3.5 | 100.0 | 100.0 |
| Forgery | 0 | 68.4 | 73.4 | 68.4 | 73.4 |
| | 1-2 | 14.9 | 9.6 | 83.2 | 83.0 |
| | 3-5 | 7.7 | 7.6 | 90.9 | 90.6 |
| | 6-10 | 4.6 | 2.3 | 95.6 | 93.0 |
| | over 10 | 4.4 | 7.0 | 100.0 | 100.0 |
| Cons | 0 | 43.5 | 44.2 | 43.5 | 44.2 |
| | 1-2 | 17.9 | 17.3 | 61.4 | 61.4 |
| | 3-5 | 11.6 | 11.1 | 73.0 | 72.5 |
| | 6-10 | 16.4 | 6.4 | 89.4 | 78.9 |
| | over 10 | 10.6 | 21.0 | 100.0 | 100.0 |
| Drug sales | 0 | 61.4 | 48.5 | 61.4 | 48.5 |
| | 1-10 | 15.4 | 9.9 | 76.4 | 58.4 |
| | 11-50 | 5.5 | 9.6 | 81.9 | 68.0 |
| | 51-100 | 6.6 | 6.1 | 88.6 | 74.1 |
| | over 100 | 11.4 | 25.9 | 100.0 | 100.0 |

N – 601–609 in 1976. N – 342–346 in 1979.

[a]The measurement period for the 1976 survey was three years; for 1979, between one and two years. All California prisoners who responded to the 1979 survey and answered these replication questions are included in this table.

Houchens, 1981, pp. 34–37). It is also consistent with research based on temporally detailed self-reports from unincarcerated populations.[6]

To further explore the validity of a model in which crimes are more likely to occur near the time of arrest, we analyzed information collected in the 1978 survey about the amount of time respondents were actually committing crimes. Respondents who indicated they had committed 11 or more crimes (of a given type) during the measurement period were then led to a question asking, "During how many of those months [in the measurement period] did you do one or more [crimes]?" (See the example for burglary in App. E.)

---

[6]See, for example, the studies of substance-abusing populations by Ball, Shaffer, and Nurco (1983), Goldstein (1982), and Johnson (1981).

Denote the answer to this question by $M_i$ for crime type i. Dividing $M_i$ by the length of the measurement period yields the approximate fraction $f_i$ time that the respondent was committing crimes of type i. Further, dividing the estimated number of crimes of type i committed during the measurement period by $M_i$ yields a rate $L_i$, which we call the individual's *peak rate* for committing crime type i. It is the rate at which he commits crime type i during the months that he commits that type of crime at all.

We then stratified respondents according to their peak rates $L_i$ and examined how their fractions $f_i$ were related to the length of their measurement period. Confirming the general appropriateness of the model described above, the data showed that among respondents in a given range of peak rate $L_i$, those who had short measurement periods were active all or nearly all of the time, while those who had long measurement periods were active smaller fractions of the time. For example, among respondents whose peak rate for drug dealing was under 360 drug deals per month,[7] those whose measurement periods were three to six months long dealt drugs on average during 80 percent of their street time, while those whose measurement periods lasted 18 to 21 months dealt drugs during 59 percent of the months on average. For business robbery, the effect of the duration of the measurement period did not vary according to the respondent's peak rate. Overall, respondents whose measurement periods were three to six months long committed business robbery on average 71 percent of the time, while those whose measurement periods were 18 to 21 months long committed business robbery 39 percent of the time.

## Adjusting Estimated Crime Rates According to the Duration of the Measurement Period

We found that a simple mathematical model adequately describes the relationship between crime commissions and length of the measurement period found in the data. In this model, the individual switches between a quiescent state and an active state for crime type i from time to time. Whenever he is in the quiescent state he has a specified probability per unit time of switching into the active state; and whenever he is in the active state he has some other probability per unit time of switching into the quiescent state.

Assuming that the end of the measurement period occurs when the offender is in the active state (because for this survey, the respondent was arrested and incarcerated for a crime committed in the last month of the measurement period), then the expected fraction of time in the active state $f_i$ is mathematically related to the length T of the measurement period by a function that declines exponentially to a steady-state value $a_i$:

$$f_i = a_i + (1 - a_i) \exp(-b_i T) \ , \tag{2.1}$$

where $a_i$ and $b_i$ are parameters that can be estimated from the data on crime type i.

When fitting this mathematical function to the data, we postulated that either $a_i$ or $b_i$ might depend on the individual's peak rate $L_i$.[8] More specifically, we allowed a slow variation with $L_i$ in the form

---

[7]This is the median peak rate for respondents whose peak rate could be calculated from the data in the survey instrument (i.e., they reported committing 11 or more drug deals during the measurement period).

[8]Strictly speaking, $L_i$ is actually the *observed* rate, while $a_i$ and $b_i$ would be envisioned to depend on the individual's *expected* peak rate. To avoid cluttering the notation, we use $L_i$ in this formulation. Note also that we are discussing a generic person rather than adding a subscript to index the individual.

$$a_i = \alpha_{1i} + \alpha_{2i} \log(L_i) \tag{2.2a}$$

$$b_i = \beta_{1i} + \beta_{2i} \log(L_i) \tag{2.2b}$$

so that in total there are four parameters in Eq. (2.1) after substitution of Eqs. (2.2). Using nonlinear least-squares fit of Eq. (2.1) with data giving $f_i$, $L_i$, and T for each person, we found that the inclusion of a log(L) term in the equation for $b_i$ did not improve the fit significantly for any crime type, whereas inclusion of the log(L) term in the form of $a_i$ improved the fit for some crime types but not others. The resulting coefficients are shown in Table 2.

For respondents whose peak rate $L_i$ for crime type i could be calculated, we defined a new *adjusted* annual crime commission rate as

$$U_i = a_i L_i \ , \tag{2.3}$$

which is an estimate of the number of crimes they would commit in a year if they were allowed to reach steady-state in their fluctuations between quiescent and active states. It differs from the value of the crime commission rate used in previous studies of RAND inmate survey data by adjusting for the length of the measurement period when estimating an individual's crime rate.

However, the format of the survey instrument (App. E) permits calculating the peak rate $L_i$ only for individuals who reported more than 10 crimes of type i committed during their measurement period. For other respondents who committed crimes of type i, we have an estimate

## Table 2

### COEFFICIENTS OF MODEL FOR ADJUSTING CRIME RATES

| Crime Type | Coefficient | | |
| --- | --- | --- | --- |
| | $\alpha_{1i}$ | $\alpha_{2i}$ | b |
| Burglary | 0.1521 | 0.0785 | 0.2128 |
| Business robbery | 0.3212 | — | 0.1325 |
| Person robbery | 0.0617 | 0.0900 | 0.3966 |
| Auto theft | 0.0573 | 0.0863 | 0.1825 |
| Other theft | 0.0371 | 0.0927 | 0.1283 |
| Forgery and credit cards | 0.1542 | 0.0355 | 0.1787 |
| Fraud | 0.3277 | — | 0.0701 |
| Drug dealing | 0.5488 | 0.0226 | 0.2182 |

NOTE: The model is:

$$f_i = \alpha_{1i} + \alpha_{2i} \log L_i + (1 - \alpha_{1i} - \alpha_{2i} \log L_i) \exp(-b_i T),$$

where $L_i$ is the peak rate for crime type i, T is the length of the measurement period, and $f_i$ is the fraction of the measurement period during which crime type i was committed. These models were fitted to RAND's 1978 survey data by (unweighted) nonlinear least squares. For business robbery and fraud, the log(L) term did not improve the fit appreciably and was not used ($\alpha_{2i} = 0$).

$N_i$ of the number of crimes they committed during a measurement period of length T. Their estimated peak rate $L_i$ is then, by definition,

$$L_i = N_i/f_iT \qquad (2.4)$$

where $f_i$ is given by Eq. (2.1) after substitution of (2.2a). Equation (2.4) has the unknown peak rate $L_i$ on both sides of the equation because it is hidden in $f_i$ on the right side. We solved the equation iteratively to obtain the value of $L_i$ (and thus $f_i$) for each respondent whose peak rate could not be determined directly from their data. The same values of the parameters $\alpha_{1i}$ and $\alpha_{2i}$ (found by fitting data for other respondents) were used in these calculations.

The adjusted crime rates of these respondents were calculated according to Eq. (2.4), just as for the respondents who reported more than 10 crimes committed. For respondents with low crime rates, this is an unverified adjustment, because data for peak rates were available only from respondents with at least moderately high crime rates. However, the present research is concerned with high-rate offenders, so errors in estimating low crime commission rates should not materially affect our results.

## Effects of Researchers' Handling of Uncertain Responses

A modest effect on the estimated distribution of annual crime commission rates and a very large influence on reported *average* estimated rates arose from researchers' handling of a small number of respondents whose data were unclear or indicated extremely high crime counts (hundreds, thousands, or even tens of thousands of crimes committed during the measurement period). For the 1976 RAND survey, the researchers established arbitrary maxima for crime rates, and any reported rates above the cutoffs were considered missing data. For the 1978 survey, the data files created from survey responses included a minimum estimate and a maximum estimate of the crime rates for each respondent and each crime type. As reported by Chaiken and Chaiken (1982, App. A), these were "not intended to be 'worst possible' cases, but rather reasonable conclusions from the data."

Analysts of the 1978 survey data used the minimum and maximum estimates in various ways. Chaiken and Chaiken (1982) sometimes used the minimum and maximum estimates separately, and sometimes they used the average of the minimum and maximum; they did not, however, specify any cutoffs above which data would be considered missing. The distributions of crime commission rates tabulated by Chaiken and Chaiken in their App. A were based on the average of the maximum and minimum, while their estimates of average crime rates are given separately for the minimum and maximum (as in Table B.1 in App. B of this report). Greenwood with Abrahamse (1982), analyzing the same data, used the maximum estimate but considered all values above the 90th percentile as missing.

In her reanalysis of the 1978 RAND survey data, Visher (1986) questioned the validity of the maximum estimate and the practice of averaging the minimum and maximum estimates together. She pointed out that the ranges from minimum to maximum were very large in some instances. For example, one California inmate had a minimum robbery rate of 82.6 and a maximum of 1,238.4, and other ranges were from 13.9 to 72.0 and from 192 to 1,032. Some wide ranges were caused by the respondent failing to answer questions fully. For example, he might check a box indicating "1 to 10" robberies committed in a street time three to six months long but not provide any answer for the next survey question asking "how many" crimes. The resulting minimum estimate is a modest $1 \times 12/6 = 2$ robberies per year, while the maximum is a seemingly prodigious $10 \times 12/3 = 40$ robberies per year.

Visher devised an alternative method for processing the 1978 survey responses; each respondent whose answer was uncertain was given a simulated answer by randomly assigning him an crime rate consistent with the distribution of answers for other respondents. She found that this method produced an answer that was much closer to Chaiken and Chaiken's minimum than to the maximum.[9]

Visher published comparative distributions of crime commission rates for inmates who reported committing robbery or burglary. Replicating the method used by Chaiken and Chaiken, she obtained a median of five robberies per year, a 90th percentile of 87 robberies, and an average falling in the range from 40.6 to 62.2 robberies. Using her method for the same respondents, she obtained a median of 3.8, a 90th percentile of 71.6, and an average of 43.4 robberies per year. Note that her average is closer to the minimum 40.6 than to the midpoint (40.6 + 62.2)/2, or 51.4.

Visher showed that researchers handling uncertain responses by calculating minimum and maximum estimates are quite likely to obtain serious overestimates of individual crime rates for the maxima but are unlikely to be far wrong with the minimum estimates.

For the purposes of classifying individuals as committing crimes at "high rates" or not at "high rates," we believe it is prudent to avoid overestimating anyone's crime rates from his data. Therefore, following the implications of Visher's work, we exclusively used the minimum estimates of crime commission rates for the remainder of this study. The minimum estimates were adjusted for the length of the measurement period as described above.

These adjusted crime rates are tabulated in App. A. They follow the same format as App. A of Chaiken and Chaiken (1982), which gives the unadjusted rates. On the whole, the effect of the adjustment is to reduce the median and 90th percentile crime commission rates (for those who commit the crime) by approximately 25 percent compared with the Chaiken and Chaiken (1982) tabulations. This change is caused primarily by our adoption of the minimum estimates, not by our adjustment for duration of measurement period (described in the previous section). For most crime types, the means of the annual commission rates reported in App. A are approximately the same as, or perhaps 5–10 percent lower than, the "minimum" estimate of the mean rate reported by Chaiken and Chaiken (1982). However, for the crimes of business robbery and fraud, the adjustment for duration of the measurement period increases the mean rate compared with the minimum estimate published in 1982.

## CLASSIFYING INDIVIDUAL CRIME RATES AS HIGH OR SERIOUS

Even after data concerning an individual's offending behavior have been processed to derive estimates of crime commission rates, questions remain concerning classifying the individuals having those rates as being "high rate" or "serious" offenders.

A threshold above which an individual's crime rate may be considered "high" can be set on grounds of both policy implications and research design. We want to classify a person's crime rate as "high" if it is well above the typical value; yet we want to avoid setting the threshold so high that the number of offenders qualifying as "high-rate" are too small to have any interest for policy purposes. For example, proposed legislation concerning a change in sentencing policy would not hold much interest if it would affect only, say, 30 or 40 convicted offenders per year in a populous state. For purposes of research design, one similarly cannot

---

[9]Because the distribution of crime commission rates is very strongly skewed to the left, the average of all responses falling in a specified range, such as from 10 crimes per year to 30 crimes per year, is much closer to the left end (10) than to the average of the two endpoints (20).

set the threshold for defining high-rate behavior so high that the sample size for "high-rate" offenders in the available data is not large enough to permit reliable statistical analysis of their characteristics.

Fortunately, the shape of the distribution of crime-commission rates permits a reasonable tradeoff between "high" and "not too high." The strongly skewed distributions characteristic of crime rates imply that most offenders who commit any particular crime type do so at fairly low rates; ordinarily the median of the distribution is substantially less than half of the 70th or 80th percentiles. Thus, individuals whose rates are near the 70th or the 80th percentile can be considered to have "high" rates compared with those near the median, which might be considered "typical."

Refer to the tabulations of medians and 75th percentiles in App. A. The 75th percentile of the burglary rate is 21.8 crimes per year, which is 4.5 times the median (4.8 burglaries per year); the 75th percentile of the business robbery rate (7.4 per year) is 3.2 times the median; and the 75th percentile for theft (51.5 crimes per year) is 9.7 times as large as the median. Similar observations could be made concerning the 70th and 80th percentiles, which are not listed in the appendix. For each of the crime types included in the survey, as many as the top 20 or 30 percent of offenders who commit the crime can reasonably be judged as committing the crime at "high" rates.

Attempting to make more specific determinations of suitable cutoff levels is unwarranted, both because there are no commonly accepted absolute standards of "high" rates and because the range from the 70th to the 80th percentile is already quite broad. For example, the range from the 70th to the 80th percentile for burglary is from 15.6 to 59.2 crimes per year, for robbery is from 8.5 to 14.9 crimes per year, and for theft is from 25.0 to 108.1. This implies, for example, that a cutoff for "high" theft rate could be set at 25, 50, 75, or 100 thefts per year, and the percentage of offenders who are classified as high-rate thieves would vary little with the cutoff that is chosen.

Our approach to classifying patterns of criminal behavior as "serious" was based on the work of Chaiken and Chaiken (1982). They found that offenders who committed robbery and assault and dealt drugs were very likely to be high-rate offenders in each of those crime types and also in other types of crimes. For this reason, they labeled the robber-assaulter-dealer a "violent predator."

We adopted six different interpretations of high-rate serious offender based on this work (Table 3). For exploratory analyses we used two different definitions of high crime rate in these interpretations: a cutoff at the 70th percentile and a cutoff at the 80th percentile. We used two different definitions of the crime rate: the unadjusted crime rate previously used by Chaiken and Chaiken (1982) and the adjusted crime rate described above and tabulated in App. A. Thus each of the six interpretations was initially calculated four different ways.

A very useful early finding from the analysis of these six interpretations was that it made little difference whether we used the adjusted or unadjusted crime rates in the definitions. To a large extent, offenders in the data who were classified as high-rate serious according to, say, interpretation 4 using the 80th percentile of the unadjusted rate as the cutoff for high rates were also classified as high-rate serious using the 80th percentile of the adjusted rate. Correspondingly few offenders had different classifications using the unadjusted 80th percentile and the adjusted 80th percentile. The relationships are shown for each of the six interpretations in Table 4.

This finding implies that the precise methods of calculating crime-commission rates are not important in defining high-rate serious offenders. The main explanation for the

16

Table 3

DEFINITIONS OF HIGH-RATE SERIOUS OFFENDERS

| Interpretation | Definition |
|---|---|
| 1 | A high robbery commission rate. |
| 2 | A high robbery rate or a high rate of robbery of persons or a high rate of robbery of businesses. |
| 3 | Interpretation 2 applies or he has a high rate of committing assault. |
| 4 | *Any* of his crime rates are high[a] and he is a violent predator (commits robbery and assault and deals drugs). |
| 5 | Any of his crime rates are high and he commits robbery and assault. |
| 6 | Any of his crime rates are high and he deals drugs and commits either robbery or assault. |

[a]That is, for any of the types of crimes studied in the RAND Second Inmate Survey.

Table 4

CLASSIFYING OFFENDERS USING ADJUSTED AND UNADJUSTED CRIME RATES[a]
(Percent of offenders whose data were suitable
for classification)

| Interpretation of High-rate Serious | Classified High-rate Serious Using 80th Percentile of | | | | Classified the Same Using Either Adjusted or Unadjusted Rates | |
|---|---|---|---|---|---|---|
| | Unadjusted Rates | | Adjusted Rates | | | |
| | N | Percent | N | Percent | N | Percent |
| 1 | 78 | 8.2 | 78 | 8.2 | 930 | 97.2 |
| 2 | 89 | 9.3 | 92 | 9.6 | 923 | 96.5 |
| 3 | 131 | 13.7 | 135 | 14.1 | 928 | 97.1 |
| 4 | 120 | 12.5 | 128 | 13.4 | 940 | 98.3 |
| 5 | 148 | 15.5 | 154 | 16.1 | 938 | 98.1 |
| 6 | 178 | 18.6 | 188 | 19.7 | 928 | 97.1 |

[a]Total sample size is 956 prisoner respondents to the RAND Second Inmate Survey for whom official record data were collected and who answered "Yes" to committing at least one of the crimes listed on the survey instrument.

insensitivity to crime-rate calculations is that the transformation from unadjusted to adjusted crime-commission rates is almost monotonic.[10] Furthermore, the domains of the distribution where the transformation is not monotonic (in fact, is far from monotonic) were found to be primarily at the very low end and the very high end. That is, the set of offenders whose rates are below the 20th percentile by one calculation of crime rates is not largely consonant with

[10]If the transformation were exactly monotonic, then by definition of percentiles the offenders whose rates are above percentile x of the unadjusted rate are also the same as the offenders whose rates are above percentile x of the adjusted rate, for any x.

the offenders below the 20th percentile by the other definition; similarly, the groups over the 95th percentile don't match up well. But we chose our high-rate cutoffs to be either at the 70th or 80th percentiles, so the adjustment rarely carries an individual's rate across the threshold. Thus, the decision to use either the unadjusted or adjusted crime-rate calculation was not very important.

These findings reinforce our decision to use only adjusted crime rates in the remainder of the study. Instead of $6 \times 2 \times 2 = 24$ interpretations of high-rate serious behavior, we are left with only $6 \times 2 = 12$ possibilities. Our strategy was to use only a few of these 12 definitions in our data analyses and to select one of the definitions as "preferred" for most of our presentations in this report. Although all of these definitions are *a priori* reasonable, whether a definition is a feasible candidate for analysis depends in part on the percentage of high-rate serious offenders in each state according to the definition. We were primarily concerned that too small a percentage would make statistical comparisons imprecise. To make our choices we tabulated these percentages by state for each of the 12 definitions; these are given in Table 5.

Table 5 shows that the definitions do, as expected, usually lead to more and more offenders being classified as high-rate serious as one moves from category 1 up through interpretation 6; interpretation 4 is the only reversal. Also, as logic implies, the above 70th percentile definitions have more offenders than the corresponding above 80th percentile definitions. As others have observed, this table indicates that California, Michigan, and Texas prison populations are successively less "hard core." The percentages in columns headed "All" refer to all offenders, and the percentages in columns headed "Reliable" refer to those

## Table 5

### PERCENTAGES OF OFFENDERS CLASSIFIED HIGH-RATE SERIOUS BY STATE

|  | | California | | Michigan | | Texas | |
|---|---|---|---|---|---|---|---|
|  | | Reliable[a] | All | Reliable | All | Reliable | All |
| 80th Percentile Interpretation[b] | | | | | | | |
|  | 1 | 16 | 15 | 8 | 9 | 3 | 3 |
|  | 2 | 17 | 16 | 9 | 10 | 4 | 4 |
|  | 3 | 25 | 24 | 14 | 15 | 6 | 6 |
|  | 4 | 23 | 21 | 12 | 14 | 8 | 7 |
|  | 5 | 27 | 26 | 15 | 17 | 9 | 8 |
|  | 6 | 29 | 28 | 23 | 23 | 14 | 11 |
| 70th Percentile Interpretation | | | | | | | |
|  | 1 | 19 | 19 | 14 | 15 | 6 | 5 |
|  | 2 | 21 | 22 | 17 | 18 | 6 | 6 |
|  | 3 | 32 | 32 | 23 | 25 | 10 | 10 |
|  | 4 | 28 | 26 | 14 | 16 | 11 | 9 |
|  | 5 | 32 | 31 | 20 | 22 | 13 | 12 |
|  | 6 | 37 | 35 | 26 | 27 | 19 | 15 |
| Sample Size | | 194 | 293 | 171 | 274 | 279 | 389 |

[a]Reliable respondents are defined by Chaiken and Chaiken (1982), App. B.
[b]See Table 3 and text for definitions of interpretations 1 through 6.

offenders whose responses were labeled "reliable" by Chaiken and Chaiken (1982).[11] No striking differences between reliable and unreliable data appear in the table.

We eliminated interpretations 1 and 2 at the 80th percentile because of the small percentage of offenders who were interpreted as high-rate serious in Texas and then selected our definitions for analysis from the remaining ten categories. For most of the exploratory analyses described in Secs. III and IV, the results were nearly the same for various different interpretations; so for this report we use four interpretations for the exploratory model specification. They are:

- Interpretation 2 using 70th percentile cutoffs of the adjusted crime rates.
- Interpretation 3 using 70th percentile cutoffs.
- Interpretation 3 using 80th percentile cutoffs.
- Interpretation 5 using 70th percentile cutoffs.

To make sure that we did not overfit the model to a particular definition of high-rate serious, we tried several different definitions. To get an independent assessment of the quality of the fit of the selected model, we used the definition based on interpretation 6 (70th percentile) in assessing our rules' performance in Secs. III and IV. In addition we computed Efron's (1986) estimate of overfit bias for our final model in Sec. IV.

---

[11]Their study compared respondents' answers on 14 topics with their official records, and checked internal consistency on 27 different items. Respondents who fell in the best 80 percent for external reliability and the best 80 percent for internal reliability were labeled "reliable." See App. B of Chaiken and Chaiken (1982) for details.

# III. OUR RESULTS IN SCREENING INDEPENDENT VARIABLES

In this section, we describe how we combined our own data analysis with previous research to specify a set of explanatory variables for use in logistic regression models to identify high-rate serious offenders. Our choice of logistic regression was a logical consequence of our decision to dichotomize offender behavior. For a broader perspective on statistical methods for predicting criminal behavior and modeling criminal careers, see Blumstein et al., Vol. I (1986), Copas and Tarling (1986), Flinn (1986), and Lehoczky (1986). Our analysis began with an examination of bivariate relationships between outcome variables and candidate explanatory variables. Explanatory variables that survived this screening were entered into a set of initial specifications of our logistic regressions. Various diagnostics were calculated from the fitted equations in order to eliminate explanatory variables that did not contribute appreciably to identifying high-rate serious offenders. This process was iterated until a satisfactory logistic regression was settled on.

We applied the model specification process to different dependent variables, to offender data from different states, and finally to different classes of explanatory variables—official record variables, self-report versions of official record variables, and the full set of self-report variables. The result is a set of logistic regression models that we use to define discrimination rules to decide which criminal offenders are high-rate serious from their explanatory variables.

## METHODOLOGY

Analysts building statistical models must consider both statistical and nonstatistical issues. In constructing discrimination rules for identifying high-rate offenders, we began with formulating the problem, deciding on the logit approach, pondering issues of acceptability of potential explanatory variables, finding acceptable proxies for explanatory variables we could not measure directly, and finally using computational and graphical methods on the data to assess how well particular models fitted the data. There are a wide variety of tools in the analyst's model-building tool box; we will briefly describe here the primary techniques that we used to screen potential explanatory or predictor variables and arrive at a set of logistic regression models that appear to fit the RAND Second Inmate Survey data.

### Examining One Explanatory Variable at a Time

Although our goal is selecting a set of explanatory variables for inclusion in a logit equation, it is frequently helpful to screen a large set of candidate variables by assessing each variable's bivariate relationship to the outcome variable being fitted. For binary outcome variables two methods are particularly useful: looking at univariate marginals of the candidate explanatory variables and examining plots of the proportion high-rate serious against each potential predictor variable.

The marginal distribution of a single explanatory variable indicates whether it has any potential as a discriminating variable. For example, a variable with 99 percent of its mass at a single value cannot possibly have very much discriminatory power. At best it can pick out 1 percent of the high-rate serious offenders. More generally, examining marginal distributions can detect errors in the data and can sometimes aid in rescaling a variable.

Before fitting ordinary regression models to continuous outcome variables, one often plots the outcome variable against the candidate predictor variable as a scatter diagram. There is no obvious analogue to a scatter plot for binary outcome variables. In order to produce a useful summary of the bivariate relationship between the two variables, some data smoothing is needed. We followed a suggestion by Copas (1983) and plotted a smoothed version of the proportion of high-rate serious offenders against each candidate explanatory variable. More formally, suppose Y(n) is equal to 1 if individual n is a high-rate serious offender and is equal to 0 otherwise. If p is the probability that Y = 1, then Copas's method plots a function of an estimate of p against a candidate explanatory variable x in order to assess the shape and strength of the relationship between p and x. In effect, this plot is a nonparametric regression of Y on x. We used Copas's smoothing method to estimate p in making these plots.

## Diagnostics for Logistic Regression

Once a set of candidate explanatory variables has been selected, the next step is to fit a logistic regression to the outcome variable of interest and assess the fit of the model to the data. In this section we describe our methods for assessing the relative strength of different possible specifications of the logistic regression equation.

For a given vector of explanatory variables x, the probability that an offender is high-rate serious is related to x by:

$$\log(p/(1 - p)) = B'x \ ,$$

where B is the vector of regression coefficients to be estimated. That is, the logarithm of the odds of being a high-rate offender is a linear function of the explanatory variables. Given a logit equation fitted by the method of maximum likelihood, we used several assessment methods. Some methods are the same as for ordinary linear regression and some are peculiar to logistic regression; we discuss them in order.

As in linear regression, computer programs for using maximum likelihood to fit a logit regression printout estimated standard errors and corresponding significance probabilities. (The analogue to the partial F-test in a linear regression is a $\chi^2$ test.) For each fitted logit equation we examined the size, sign, and statistical significance of each estimated coefficient. We compared coefficient estimates in successively nested equations looking for anomalies and evidence of colinearity problems. We discarded candidate predictor variables whose estimated coefficients were small and not statistically significant. Exceptions to this practice were variables of special interest (e.g. some of the official record variables) and variables of policy interest. In these cases we kept the variable regardless of coefficient size. We describe below how models specified with one data set were validated on another.

The second assessment step consisted of examining individual data points. In linear regression, this is often done with residual plots and looking for extreme data points—either in the explanatory variable space (high leverage points) or in outcome space (outlying residuals). There are well understood and accepted diagnostic methods for linear regression that are described in detail in several texts (for example, Belsley, Kuh, and Welch, 1979; Cook and Weisberg, 1982; and Weisberg, 1985). However, the tools currently available for logit models are not as readily accessible or as informative as those for linear regression. See Landwehr, Pregibon, and Shoemaker (1984), Pregibon (1981), Reboussin (1984), and Rosenbaum and Rubin (1983b) for recently developed logistic regression diagnostic methods. For diagnostics, we elected to use propensity plots as suggested by Reboussin (1984) and Rosenbaum and Rubin (1983b).

The value of propensity plots is based on a theorem of Rosenbaum and Rubin (1983a). It states that for a fixed propensity score (Prob($Y = 1 \mid x$)) the two (multivariate) distributions of the vector of explanatory variables corresponding to $Y = 1$ and $Y = 0$ should be the same. It therefore follows that for the correct model, a plot of the estimated propensity score (the logit probability) on the horizontal axis against a particular variable on the vertical axis that distinguishes between cases where $Y = 0$ and $Y = 1$ should have the two distributions of points on each vertical line that are similar. Dissimilarities between the two sets of values of explanatory variables corresponding to outcome values 0 and 1 for particular values of the propensity score indicate what range of probabilities the logistic model fits poorly. We plotted the estimated propensity score against each explanatory variable as part of our model specification process. To give an overall assessment of the quality of fit, we plotted the estimated propensity score against a smoothed version of the outcome variable. We used propensity plots as our major diagnostic tool in assessing successive logit fits to our data.

## SCREENING OF CANDIDATE VARIABLES

This section describes our analytic approach in regard to the following issues:

- Whether to exclude possibly unreliable observations,
- Which definition of a high-rate serious offender to use at the exploratory and final stages of fitting the logit models,
- Which states to fit initially, and
- Whether to allow full interactions in the logit models between different sets of explanatory variables and states.

### Principles of Variable Selection

In Sec. II we described how we arrived at our four primary definitions of high-rate offenders primarily on the basis of the proportions of high-rate offenders in each of our three study states. We made three additional decisions in structuring the screening of explanatory variables.

First, measures of internal and external reliability of RAND's Second Inmate Survey had previously been calculated by Chaiken and Chaiken (1982, App. B). All model screening was done after removing respondents whose data were judged to have poor internal or external reliability in the earlier analysis. Our rationale was that we did not want choices of explanatory variables to be determined by possible anomalies in the unreliable data. When we had settled on a final set of models, we refitted them to all the data, both reliable and unreliable. We observed no appreciable differences between fitted logit models from the two datasets. Therefore, to give a larger sample size we elected to use the full dataset for the results reported in this and the next section.

Second, explanatory variables with strong discriminating power in any one of the three states were included in the models for each of the three states. Using the same set of explanatory variables across states makes comparisons meaningful. One result of this uniformity constraint is that our logit equations have somewhat more variables in them than would have been the case if each state's data had separately determined the explanatory variables. Because of the differences in the three states' legal and administrative practices with respect to the criminal justice system, we decided to estimate separate logistic regression models for each state

rather than use, say, a dummy variable for each state. That would assume that the probability of being a high-rate serious offender varies with offender characteristics in the same way in all states. In statistical terms, we allowed full interactions with all variables in the model.

Third, in structuring the screening process for independent variables we fitted models with three classes of independent variables: official record variables, self-report versions of official record variables, and a full set of self-report variables. Our results in selecting explanatory variables are organized by these three classes below.

Finally, we first used the 70th and 80th percentile of interpretation 3 (high rate of robbery, robbery of persons, robbery of businesses, or committing assault) in fitting candidate logit models. We then validated the models by fitting the outcome variable of the 70th percentile of interpretation 6 (high rate of any crime plus drug dealing and robbery or assault). We focus primarily on results for this last outcome variable in our description below.

## Structure and Results

We present the results of our selection of variables for the logit models in three categories: official record variables, self-report versions of official record variables, and a full self-report set of variables. For each category of variables we prepared a list of candidate variables based on previous research on identifying high rate offenders as discussed in Sec. II. We gave substantial weight to the variables that showed any discriminatory power in predicting crime commission rates in Chaiken and Chaiken (1982). Our goal was to eliminate variables with low discriminatory power and where possible combine low information content variables into a single variable with modest or even high discriminatory power. For example we looked for combinations of arrests for a variety of different crimes to produce a single crime index variable. We used California data first in our model specification process. This section concludes with a summary of comparisons of the logit models across states and across dependent variables.

*Official Record Explanatory Variables.* Examining marginals and the Copas "p vs. x" plots for each candidate explanatory variable led us to eliminate many variables and to combine and rescale others. Table 6 presents a summary of how the official record variables discriminate among high-rate offenders in California, Texas, and Michigan and others for our primary outcome variable (70th percentile interpretation 6: high rate for any crime plus dealing drugs and robbery or assault).

The entries in Table 6 give the estimated factor that a one unit increase in the predictor variable will change the odds of the inmate being identified as high-rate serious by the logistic regression equation. For example, because 26 percent of the Michigan inmates are high-rate serious by this definition, a typical Michigan inmate has odds of .35(=.26/.74) to 1 of being high-rate serious. All other things equal, each additional past conviction would increase the odds of a Michigan inmate being high-rate serious by 22 percent. For the typical Michigan inmate, this one additional conviction changes the estimated odds from .35 to .43(1.22 × .35) and the corresponding probability from .26 to .30(= .43/1.43).

Table 6 reveals that official record information has a modest but statistically significant relationship with being a high-rate serious offender in each of the three states.

For explanatory variables, we selected four justice system variables that had the strongest predictive power from among a wide variety of candidate variables in our exploratory analysis. As the table shows, the four selected predictor variables are based on number of past convictions, time between first recorded arrest and the reference period, and total robbery and assault

## Table 6

### MULTIPLICATIVE ODDS FOR LOGIT REGRESSION USING OFFICIAL RECORD VARIABLES

(Outcome variable: 70th percentile interpretation 6 of
high-rate serious)

| Variable | California | Michigan | Texas |
|---|---|---|---|
| Total past convictions[a] | 1.00 | 1.22 | 1.04 |
| | ( )[b] | (+) | ( ) |
| Square of (robbery arrests during reference period) | 1.07 | 5.63 | H[d] |
| | (+) | ( ) | (++) |
| Square of (assault arrests during reference period) | 5.37 | .82 | 1348. |
| | (+) | ( ) | ( ) |
| Years since first recorded arrest[c] | 1.01 | .95 | 1.13 |
| | (+) | ( ) | (+) |
| Reference period ≤ 4 months? | 1.01 | 1.44 | 1.04 |
| | (+) | (+) | ( ) |
| Age during the reference period | .98 | .91 | .83 |
| | (− −) | (−) | (− −) |
| Squared age | 1.00 | .99 | 1.00 |
| | ( ) | ( ) | ( ) |
| Intercept | 2.12 | 3.66 | 7.26 |
| | (+++) | (+++) | (+) |
| Percent high rate | 36% | 26% | 15% |
| (odds) | (.56/1) | (.35/1) | (.17/1) |
| Sample size | 275 | 244 | 357 |
| $\chi^2$ (7 d.f.)[e] | 16.6 | 17.8 | 20.0 |

[a]For assault, burglary, drugs, murder, rape, robbery, or kidnap.

[b]The entries within the parentheses correspond to the levels of the t-statistics: +++ = t > 3.0; ++ = 3.0 ≥ t > 2.0; + = 2.0 ≥ t > 1; and blank is 1 or less. The minuses are defined correspondingly.

[c]Years between first recorded arrest and beginning of reference period.

[d]"H" is used for any coefficient estimate above 10,000. Such high estimates occur when only a single case or two have the value 1 and also have an extreme value of the independent variable.

[e]The percentage points for a $\chi^2$ distribution with 7 degrees of freedom are 14.0, 16.0, and 18.5 for significance levels of .05, .025, and .01 respectively.

arrests during the reference period. We also tried several potential explanatory variables based on arrests and convictions for other crime types, but they yielded poorer fits. As measured by statistical significance and overall contribution to fit, the justice system variables contribute mildly to identifying high-rate offenders. However, for some of the variables, the presence of a high value changes the estimated probability substantially (see the Michigan example using past convictions above). This phenomenon occurs because few inmates have high values of these variables and hence estimates for most inmates are unaffected by these variables.

The "measurement" variable (whether the inmate's reference period was short) and two age variables round out the explanatory variables in the equation. All else equal, inmates with shorter street times are more likely to be high-rate serious offenders. Shorter street times are also indicators of recent incarceration experiences, which may partially explain this result. To the extent that street time is measuring time free (and its complement—time served during the last two years), this variable could be regarded as an official record variable. Also, all else equal, younger inmates are more likely to fall into the high-rate serious category than older ones.

The equations for the three states exhibit no striking differences in coefficient patterns.

The same explanatory variables were used to estimate logistic regressions for two of the other more stringent definitions (70th and 80th percentile of interpretation 3) of high-rate serious offenders in the three states. The coefficient patterns were similar. The age variables contributed most to discriminating between those in the high-rate category and those who were not. In most of the logistic regressions, having a street time of four months or less gave good discriminatory power—in one case changing the odds by a factor of five. Overall, the finding of only modest discriminatory strength based on official record data (Chaiken and Chaiken, 1982) also holds for the other definitions of high-rate serious offenders.

*Self-Report Official Record Explanatory Variables.* In their analysis of the RAND Inmate Survey data, Chaiken and Chaiken (1982) found that a prisoner's self-reports of official record information was more indicative of his self-reported crime commissions than the official records themselves. To the degree that self-report versions of official record variables yield improved identification of high-rate serious offenders, this raises the possibility that better record keeping might increase discriminatory power.

Table 7 shows that our logit regression results are qualitatively similar to what Chaiken and Chaiken got using linear regression with crime rate as their outcome variable. That is, comparing the fitted coefficients for the self-report variables in Table 7 with the corresponding official record variables in Table 6 reveals generally stronger discriminatory power for the self-report versions. Robbery arrests, assault arrests, and length of criminal career (years since first arrest) all have estimated coefficients that are of greater magnitude and of greater statistical significance in their self-report versions than the official record versions. (Robbery arrests are represented by both a linear and squared term in their self-report versions and by only a squared term in the official record version because a linear term contributed nothing to the fit.) Consistent with the official record logit regression modeling, other combinations of arrests and convictions for various crime types did not contribute very much to the explanatory power of the equations. As was true in the official record model, age and whether the reference period was short were moderately strong predictor variables. The better fit resulting from using the self-report versions of the justice system variables may reflect the fact that the outcome variable is also based on self-reported crime commissions.

*Unrestricted Explanatory Variable Set.* Based largely on Chaiken and Chaiken (1982) we started with a substantial list of candidate explanatory variables in developing the full self-report model. Table 8 shows that beyond self-report official record information, nine additional explanatory variables describe each offender's background and we use them to identify high-rate serious offenders. Three of the selected variables capture aspects of the inmate's earlier criminal career—whether he committed violent crime frequently as a juvenile, whether he was committed to a state juvenile facility, and the number of months during the two years before the reference period that he was incarcerated. Two variables relate to substance abuse—whether he was a frequent heroin user as a juvenile and whether he abused barbiturates and alcohol daily during the reference period. Other background variables include whether he graduated from high school, the percentage of the reference period he was employed, and whether he is black or hispanic. Minority or ethnic status contributes to a slightly *lower* risk that an inmate is a high-rate serious offender. As shown in Sec. IV, these nine variables add substantially to the discriminatory power of the identification rule based on the logistic regression equations.

## Table 7

### MULTIPLICATIVE ODDS FOR LOGIT REGRESSION USING SELF-REPORTED OFFICIAL RECORD VARIABLES

(Outcome variable: 70th percentile interpretation 6 of
high-rate serious)

| Variable | California | Michigan | Texas |
|---|---|---|---|
| Number of prior felony convictions | 1.10 | 1.07 | 1.00 |
| | (+)[a] | ( ) | ( ) |
| Self-reported robbery arrests during | 89.77 | 428 | H[c] |
| reference period | (+) | (++) | (++) |
| Square of self-reported robbery arrests | H | .11 | .99 |
| during reference period | (−) | (−) | (− −) |
| Square of self-reported assault arrests | H | H | 176 |
| during reference period | (+) | (+) | ( ) |
| Years since first self-reported arrest[b] | 1.09 | 5.58 | 1.13 |
| | (+) | (++) | (++) |
| Reference period ≥ 4 months? | 1.52 | 1.92 | 1.34 |
| | (+) | ( ) | ( ) |
| Age during the reference period | .88 | .77 | .83 |
| | (− −) | (− − −) | (− − −) |
| Aged squared | .99 | .98 | 1.00 |
| | (−) | (− −) | ( ) |
| Intercept | 4.47 | 29.70 | 4.62 |
| | (+) | (++) | (+) |
| Percent high rate | 36% | 26% | 15% |
| (odds) | (.56/1) | (.35/1) | (.17/1) |
| Sample size | 275 | 244 | 357 |
| $\chi^2$ (8 d.f.)[d] | 24.1 | 19.2 | 25.8 |

[a]The entries within the parentheses correspond to the levels of the t-statistics: +++ = t > 3.0; ++ = 3.0 ≥ t > 2.0; + = 2.0 ≥ t > 1; and blank is 1 or less. The minuses are defined correspondingly.

[b]Years between first self-reported arrest and reference period.

[c]"H" is used for any coefficient estimate above 10,000. Such high estimates occur when only a single case or two have the value 1 and also have an extreme value of the independent variable.

[d]The percentage points for a $\chi^2$ distribution with 8 degrees of freedom are 15.5, 17.5, and 20.1 for significance levels of .05, .025, and .01 respectively.

## Table 8

### MULTIPLICATIVE ODDS FOR LOGIT REGRESSION USING ALL SELF-REPORTED VARIABLES

(Outcome variable: 70th percentile interpretation 6 of high-rate serious)

| Variable | California | Michigan | Texas |
|---|---|---|---|
| Number of prior felony convictions | .97 | 1.15 | .96 |
| | ( )[a] | ( ) | ( ) |
| Age during the reference period | .87 | .94 | .85 |
| | (− −) | (−) | (− −) |
| Aged squared | .99 | .99 | 1.00 |
| | ( ) | ( ) | ( ) |
| Years since first self-reported arrest[b] | 1.08 | .98 | 1.10 |
| | (+) | ( ) | (+) |
| Reference period < 4 months? | 1.27 | 1.12 | 1.66 |
| | ( ) | ( ) | ( ) |
| Self-reported robbery arrests during reference period | 66.02 | 14.16 | H[c] |
| | (+) | ( ) | (++) |
| Square of self-reported robbery arrests during reference period | .00 | .76 | .98 |
| | ( ) | ( ) | (− −) |
| Square of self-reported assault arrests during reference period | H | .01 | 6.64 |
| | (+) | ( ) | ( ) |
| Violent crime frequently as juvenile | 2.85 | 3.64 | 1.50 |
| | (++) | (++) | ( ) |
| Commitment to state juvenile facility | .53 | 1.12 | 1.22 |
| | (−) | ( ) | ( ) |
| Number of months incarcerated in window 1[d] | 1.12 | .87 | .91 |
| | (+) | ( ) | ( ) |
| Juvenile frequent heroin use? | 2.00 | 1.06 | 3.15 |
| | (+) | ( ) | (++) |
| Graduated high school? | 1.45 | 1.81 | 1.34 |
| | (+) | (+) | ( ) |
| Percent of time employed during reference period | .22 | 1.07 | .39 |
| | (− − −) | ( ) | (−) |
| Barbiturates and alcohol abuse daily during reference period | 2.43 | 727 | 6.39 |
| | (+) | ( ) | (+++) |
| Black? | .29 | .43 | .35 |
| | (− − −) | (− −) | (− −) |
| Hispanic? | .67 | .40 | 1.14 |
| | (−) | ( ) | ( ) |
| Intercept | 10.48 | 1.54 | 5.45 |
| | (++) | ( ) | (+) |
| Percent high rate | 36% | 26% | 15% |
| (odds) | (.56/1) | (.35/1) | (.17/1) |
| Sample size | 275 | 244 | 357 |
| $\chi^2$ value (17 d.f.)[e] | 33.3 | 52.4 | 65.8 |

[a]The entries within the parentheses correspond to the levels of the t-statistics: +++ = t > 3.0; ++ = 3.0 ≥ t > 2.0; + = 2.0 ≥ t > 1; and blank is 1 or less. The minuses are defined correspondingly.

[b]Years between first self-reported arrest and reference period.

[c]"H" is used for any coefficient estimate above 10,000. Such high estimates occur when only a single case or two have the value 1 and also have a high value of the independent variable.

[d]A two year period ending two years before the reference (also known as window 3) period.

[e]The percentage points for a $\chi^2$ distribution with 17 degrees of freedom are 27.6, 30.2, and 33.4 for significance levels of .05, .025, and .01 respectively.

# IV. USING FITTED MODELS TO IDENTIFY HIGH-RATE SERIOUS OFFENDERS

The fitted logit regression model yields estimated probabilities that offenders with given attributes will be high-rate serious. Although the estimated probabilities are of some interest in their own right, our primary interest here is the performance of discriminant rules based on the estimated probabilities. Therefore, we define threshold values of the probabilities that yield discriminant rules—e.g., all offenders whose estimated probabilities of being high-rate serious is above 0.50 might be labeled as high-rate serious offenders.

For each of the main logit regressions of interest, we present contingency tables giving the discrimination results of the rules based on the equations and give a measure of the discriminatory power of the rule on the data.

## RESULTS FOR OUR PREFERRED DEFINITION OF HIGH-RATE SERIOUS

We present results from models of the 70th percentile interpretation 6 definition of high-rate serious (high rate of any crime plus drug dealing and robbery or assault) using three different sets of predictor variables: official record variables, self-report versions of official record variables, and the full self-report set of variables. Comparing offender identifications and corresponding measures of discriminatory power across the three sets of variables allows us to quantify the additional value of the self-report information compared with the official record data on offenders.

The logistic regression equations and the propensity score plots based on the fitted equations are necessary ingredients for defining a discrimination rule. We define this rule by setting a threshold value; offenders whose fitted probabilities are above the threshold are identified as high-rate serious offenders, and offenders whose fitted values are below the thresholds are labeled as nonhigh-rate serious offenders. For each discriminant rule, we also compute measures of discriminatory accuracy and compare these measures across states and across sets of predictor variables. In this section, we compare the accuracy of rules based on official record predictors, self-report official record predictors, and a full set of self-report predictor variables.

### Rule Definitions and Measures of Accuracy

In defining a discriminant rule, there is a tradeoff between false positives and false negatives. For example, in Texas only 15 percent (53) of the 357 offenders were high-rate serious offenders using our preferred definition. If reducing false positives is paramount, a discriminant rule with only a 15 percent error rate can be constructed by labeling all offenders as not being high-rate serious. The resulting table is:

|               |      | Actual |       |
| ------------- | ---- | ------ | ----- |
|               |      | High   |       |
| Predicted     | Not  | Serious |      |
| Not           | 304  | 53     | 357   |
| High Serious  | 0    | 0      | 0     |
| Totals        | 304  | 53     | 357   |

We will adopt a common convention to resolve this tradeoff; we constrain our discriminant rules to identifying the same proportion high-rate serious as are actually high-rate serious—i.e., matching the selection rate to the base rate. The middle panel of Table 9 gives the resulting contingency table for the Texas data using a constrained rule with the official record logistic regression equation (Model 1).

Table 9 uses the logistic regressions given in Table 6. The "threshold" matches the selection rate to the base rate. It gives the estimated probability (propensity score) that separates the offenders identified as high-rate serious from those identified as not high-rate serious in each state. Because there are different mixes of offenders with respect to seriousness of crimes in prison in California, Texas, and Michigan, the threshold values vary. Inmates are classified as high-rate serious if their predicted p ≥ threshold. A separate threshold is chosen in each state so that the fraction of the sample predicted to be high-rate serious is the same as the actual fraction of high-rate serious offenders in that state's sample. Indeed, we see that Texas has the most stringent threshold of .20 compared with .39 in California and .34 in Michigan. The more stringent threshold for Texas reflects the high percentage of less serious offenders sentenced to prison in Texas—i.e., offenders who might get probation in another jurisdiction tend to go to prison in Texas.

For 2×2 contingency tables such as those given in Table 9, there are many definitions of predictive accuracy in the statistical literature (Fienberg, 1977). The criminal justice literature on recidivism and predicting high-rate serious offenders mainly uses the "Relative Improvement Over Chance" or RIOC (Fischer, 1984b; and Gottfredson and Gottfredson, 1986). The RIOC can be interpreted as the proportional improvement over chance in discriminatory

## Table 9

PREFERRED DEFINITION OF HIGH-RATE SERIOUS OFFENDER:
MODEL 1 IDENTIFICATIONS

| | California | | Texas | | Michigan | | All States | |
|---|---|---|---|---|---|---|---|---|
| | Actual | | Actual | | Actual | | Actual | |
| Predicted | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Not | 123 (70) | 53 (54) | 269 (89) | 35 (66) | 139 (77) | 41 (64) | 531 (80) | 129 (60) |
| High Serious | 53 (30) | 46 (46) | 35 (12) | 18 (34) | 41 (23) | 23 (36) | 129 (20) | 87 (40) |
| Total | 176 (64) | 99 (36) | 304 (94) | 53 (6) | 180 (85) | 64 (15) | 660 (75) | 216 (25) |
| Threshold | .39 | | .20 | | .34 | | — | |
| RIOC[a] | .16 | | .22 | | .13 | | .21 | |
| Sample size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for definition of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.
[a]RIOC = Relative improvement over chance.

power; its minimum value of 0.00 indicates no improvement, and its maximum value of 1.00 is perfect discrimination. For a dichotomous outcome as we are using, the definition is

$$\text{RIOC} = \frac{P - P_c}{1 - P_c}$$

where P is the proportion of cases correctly identified and $P_c$ is the proportion of cases correctly labeled by a chance rule. For example, the middle panel of Table 9 giving the predictors based on Texas official record data has an RIOC given by

$$\text{RIOC} = \frac{\dfrac{269+18}{357} - \left(\dfrac{304}{357}\right)^2 - \left(\dfrac{53}{357}\right)^2}{1 - \left(\dfrac{304}{357}\right)^2 - \left(\dfrac{53}{357}\right)^2}$$

$$= \frac{.804 - .747}{1 - .747}$$

$$= .2247$$

Note that the actual proportion labeled is $P = (269 + 18)/357 = 0.80$ while the proportion correctly labeled by a chance rule is defined as the 85 percent of the 85 percent low-rate offenders plus 15 percent of the 15 percent high-rate offenders for a total of 74.7 percent. The intuitive appeal of RIOC is its analogy with $R^2$ or the coefficient of determination as a measure of goodness of fit in linear regression.

There are various other measures of association in the statistical literature for 2×2 contingency tables. Most of these are functions of either the familiar $\chi^2$ statistic or the cross product odds ratio. See Bishop, Fienberg, and Holland (1975, Chapter 11) for a discussion.

## Discriminant Rule Results

*Results for Official Record Variables.* Table 9 shows the results of classifying offenders using official record variables. It reveals patterns consistent with the earlier Chaiken and Chaiken (1982) and Greenwood and Abrahamse (1982) analyses of the RAND Second Inmate Survey data; discrimination rules of this type make many fewer mistakes for the low-rate offenders than for the high-rate offenders. The relative error rates for the two groups are 30 to 54 percent in California, 12 to 66 percent in Texas, and 23 to 64 percent in Michigan. In short, discrimination rules based on official record information give high error rates (60 percent overall) in identifying high-rate serious offenders compared with a modest 20 percent overall error rate in correctly labeling nonhigh-rate serious offenders. The combined RIOC measure of 21 percent quantifies the low discriminant power that official records have for identifying high-rate serious offenders.

*Results for Self-Report Official Record Variables.* The classification rules based on the self-report versions of official records again agree with earlier analyses of the RAND Second Inmate Survey data. Again, the rules make mistakes much less often for nonhigh-rate serious offenders than for high-rate serious offenders—29 to 52 percent in California, 11 to 64 percent in Texas, and 21 to 59 percent in Michigan. Overall, Table 10 shows that 19 percent of low rate offenders are incorrectly labeled, and a whopping 57 percent of the high-rate serious offenders are mislabeled. This results in an overall RIOC of 25 percent.

Table 10

| | California | | Texas | | Michigan | | All States | |
|---|---|---|---|---|---|---|---|---|
| | Actual | | Actual | | Actual | | Actual | |
| Predicted | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Not | 125 | 51 | 270 | 34 | 142 | 38 | 537 | 123 |
| | (71) | (52) | (89) | (64) | (79) | (59) | (81) | (57) |
| High Serious | 51 | 48 | 34 | 19 | 38 | 26 | 123 | 93 |
| | (29) | (49) | (11) | (36) | (21) | (41) | (19) | (43) |
| Total | 176 | 99 | 304 | 53 | 180 | 64 | 660 | 216 |
| | (64) | (36) | (85) | (15) | (74) | (26) | (75) | (25) |
| Threshold | .40 | | .23 | | .33 | | — | |
| RIOC | .20 | | .25 | | .20 | | .25 | |
| Sample size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for interpretation of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.

In aggregate, Model 2 gives a modest improvement in discriminant rule accuracy. Of the 876 offenders, Table 10 shows that six corresponding erroneous classifications of low rate and six erroneous classifications of high-rate serious offenders are corrected compared with Model 1 (the official record) identifications. The RIOC correspondingly increases from 21 to 25 percent with the use of the self-report versions of the official record variables.

*Results from the Full Self-Report Variables Model.* Table 11 shows that the Model 3 (full self-report logistic regression) rules improve on using only official record information. Again, the batting averages are higher in correctly identifying nonhigh-rate serious offenders than high-rate serious offenders. The state with the most notable improvement is Michigan where *only* 45 percent of the high-rate serious offenders and 16 percent of the nonhigh-rate serious offenders were mislabeled. (The corresponding Model 2 rule numbers for Michigan are 59 percent and 21 percent.) Overall, the RIOC value is a more respectable 39 percent for rules using the full self-report set of variables—a substantial increase from the 25 percent value for the Model 2 (self-report official variables) rule.

Comparing Tables 12 and 11 confirms our earlier conclusions that knowing the offender's race makes little, if any, contribution to correctly labeling him beyond the weak discriminant rule we have been able to construct. Overall, the inclusion of race actually leads to more incorrect identifications, 208 rather than 200. (Although the RIOC in Michigan decreases from .39 to .30, one should not make too much of a change in the classification of eight inmates. The decrease when using an additional explanatory variable is an artifact of changes in the particular threshold value used for the discriminant rule.) It is reassuring that the set of predictors in the full self-report model is sufficiently rich in describing the offenders that the respondent's race is not needed to improve the performance of the discriminant rule.

## Table 11

### PREFERRED DEFINITION OF HIGH-RATE SERIOUS OFFENDER: MODEL 3 IDENTIFICATIONS

| Predicted | California Actual | | Texas Actual | | Michigan Actual | | All States Actual | |
|---|---|---|---|---|---|---|---|---|
| | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Not | 135 (77) | 41 (41) | 274 (90) | 30 (57) | 151 (84) | 29 (45) | 560 (85) | 100 (46) |
| High Serious | 41 (23) | 58 (59) | 30 (10) | 23 (43) | 29 (16) | 35 (55) | 100 (15) | 116 (54) |
| Total | 176 (64) | 99 (36) | 304 (85) | 53 (15) | 180 (74) | 64 (26) | 660 (75) | 216 (25) |
| Threshold | .39 | | .25 | | .32 | | — | |
| RIOC | .35 | | ,34 | | .39 | | .39 | |
| Sample Size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for interpretation of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.

## Table 12

### PREFERRED DEFINITION OF HIGH-RATE SERIOUS OFFENDER: MODEL 4 IDENTIFICATIONS

| Predicted | California Actual | | Texas Actual | | Michigan Actual | | All States Actual | |
|---|---|---|---|---|---|---|---|---|
| | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Not | 135 (77) | 41 (41) | 274 (90) | 30 (57) | 147 (82) | 33 (52) | 556 (84) | 104 (48) |
| High Serious | 41 (23) | 58 (59) | 30 (10) | 23 (43) | 33 (18) | 31 (48) | 104 (16) | 112 (52) |
| Total | 176 (64) | 99 (36) | 304 (85) | 53 (15) | 180 (74) | 64 (26) | 660 (75) | 216 (25) |
| Threshold | .42 | | .25 | | .31 | | — | |
| RIOC | .35 | | .34 | | .30 | | .36 | |
| Sample size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for interpretation of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.

*Overall Performance Comparisons of Rules.* Figure 1 summarizes the performance of the three discrimination rules. Because race variables do not improve performance as described above, we do not include the rule corresponding to the full self-report plus race model in these comparisons. The RIOC for Model 1, the official record logistic regression, is 21 percent—about a fifth of the way up from coin tossing to perfect prediction. Model 2, the self-report official record version (chosen as possibly indicative of the performance of a rule with better quality official record data) has an RIOC of 25 percent, a four percentage point improvement over Model 1. The RIOC for the full self-report rule (Model 3) can be thought of as an upper bound on how well the criminal justice system might do with fairly complete information on offenders. How close one could get to the upper bound in practice depends on legal and ethical considerations relating to what offender information can be used in sentencing. The RIOC for Model 3 is 39 percent, almost double that for Model 1. Thus, although improving on official record information clearly has the potential for substantially improved predictions of who is a high-rate serious offender, the most one can hope for is getting less than half the way from "coin tossing" to perfect prediction. That is, even the Model 3 (full self-report) rule incorrectly labels 23 percent (200) of the offenders, compared with chance where 36 percent are incorrectly labeled. This is not an encouraging track record.

As one might expect, the improvement of the Model 3 rule over the Model 1 rule is greatest in absolute terms for high-rate serious offenders compared with nonhigh-rate serious offenders. That is, the 40 percent rises to 54 percent of the high-rate serious offenders that are correctly labeled, and the 80 percent correctly identified nonhigh-rate serious offenders rises to 85 percent under Model 3. As Fig. 1 shows, all three discrimination rules perform far better with nonhigh-rate serious offenders than high-rate serious offenders.



Fig. 1—Performance of three discrimination rules using preferred definition of high-rate serious offenders

## SENSITIVITY ANALYSIS AND CROSS VALIDATION

Two potential problems may arise in prediction studies. First, how sensitive is the accuracy of predictions to the precise definition of the outcome variable (high-rate serious offender)? As reported in Sec. II, we experimented with 12 different definitions. While details differ, these rules perform similarly to the rule based on our preferred definition of high-rate serious offender. Here, we present results for one other definition to illustrate this point. A second issue we will briefly discuss is cross validation. That is, how well would our measures of performance of the classification rule apply to a new independent set of offenders?

### Other Definitions of High-Rate Serious Offender

Table 13 gives the results of using the 80th percentile of interpretation 3 definition of high-rate serious offender (see Table 3) with the official record predictor variables used in Model 1. That is, an offender is defined as high-rate serious if he has a high robbery rate or a high rate of robbery of persons or a high rate of robbery of businesses or a high rate of committing assaults. This is more stringent than our preferred definition and yields 14 percent of the offenders being high-rate serious (rather than 25 percent with the preferred definition).

In terms of improvement over the "chance rule," the overall RIOC of 29 percent for the alternative definition rule is higher than the RIOC of the preferred definition rule of 21 percent. However, this comparison is difficult to interpret, because the RIOC is sensitive to variations in the selection and base rates used even on the same data. Indeed, RIOC values generally improve as the selection rate gets more extreme.

### Table 13

ALTERNATIVE DEFINITION OF HIGH-RATE SERIOUS OFFENDER:
MODEL 1 IDENTIFICATIONS

| | California | | Texas | | Michigan | | All States | |
|---|---|---|---|---|---|---|---|---|
| | Actual | | Actual | | Actual | | Actual | |
| Predicted | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Not | 174 (84) | 34 (51) | 321 (95) | 16 (80) | 183 (88) | 25 (69) | 678 (90) | 75 (61) |
| High Serious | 34 (16) | 33 (49) | 16 (5) | 4 (20) | 25 (12) | 11 (31) | 75 (10) | 48 (39) |
| Total | 208 (76) | 67 (24) | 337 (94) | 20 (6) | 208 (85) | 36 (15) | 753 (86) | 123 (14) |
| Threshold | .30 | | .13 | | .22 | | — | |
| RIOC | .33 | | .15 | | .19 | | .29 | |
| Sample size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for interpretation of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.

Table 14 shows the results of using the full self-report variables (Model 3) with this definition of high-rate serious offender. The Model 3 rule performance improvement over Model 1 is similar for both definitions of high-rate serious offenders; the overall RIOC here increases from 29 percent to 49 percent. Also similarly, most of the improvement is in reducing misclassifications of high-rate serious as opposed to nonhigh-rate serious offenders. In sum, the results for the rule based on this alternative definition of high-rate serious hold no surprises.

## Cross Validation

How well would our measures of performance of the classification rule apply to a new independent set of offenders? That is, the apparent error rate in classifying the offenders that were used to derive the rule will usually underestimate the true error rate that would obtain in applying the rule to a different set of offenders. One crude check on this question is comparing how consistent our logit equations are for the different outcome variables. Using one definition of the high-rate serious offender to select the explanatory variables and another to "calibrate" the equation gives us some protection.

In any potential application of our rules for sentencing policy, obviously offenders other than RAND Survey respondents would be used. Efron (1986) has derived a simple estimate of the downward bias in the apparent error rate in such situations, for discrimination rules based on logistic regressions. Efron's theoretical results and his simulation results indicate that this error rate estimate is quite accurate and so can be used with confidence. Table 15 gives the bias estimates in the three states for the preferred definition of high-rate offenders. The bias estimates for the alternative definition of high-rate serious offender (described in Tables 13 and 14) are of the same magnitude as those shown in Table 15.

Table 14

ALTERNATIVE DEFINITION OF HIGH-RATE SERIOUS OFFENDER:
MODEL 3 IDENTIFICATIONS

| | California | | Texas | | Michigan | | All States | |
|---|---|---|---|---|---|---|---|---|
| | Actual | | Actual | | Actual | | Actual | |
| Predicted | Not | High Serious | Not | High Serious | Not | High Serious | Not | High Serious |
| Low | 184 (87) | 27 (40) | 327 (97) | 10 (50) | 191 (92) | 17 (47) | 699 (93) | 54 (44) |
| High | 27 (13) | 40 (60) | 10 (3) | 10 (50) | 17 (8) | 19 (53) | 54 (7) | 69 (56) |
| Total | 208 (76) | 67 (24) | 337 (94) | 20 (6) | 208 (85) | 36 (15) | 753 (86) | 123 (14) |
| Threshold | .38 | | .10 | | .30 | | — | |
| RIOC | .47 | | .47 | | .45 | | .49 | |
| Sample size | 275 | | 357 | | 244 | | 876 | |

NOTE: See Table 3 for interpretation of high-rate serious. Numbers in parentheses are the relative percentages of the low or high-rate offenders labeled low or high.

## Table 15

BIAS ESTIMATES OF ERROR RATE FOR PREFERRED
DEFINITION OF HIGH-RATE SERIOUS OFFENDER

| Model | California | Texas | Michigan |
|---|---|---|---|
| Model 1: Official Records | .0344 (.3855)[a] [42%][b] | .0206 (.1961) [22%] | .0349 (.3361) [37%] |
| Model 2: Self-report Official Records | .0325 (.3709) [40%] | .0170 (.1905) [21%] | .0357 (.3115) [35%] |
| Model 3: Full self-report | .0404 (.2982) [34%] | .0220 (.1681) [19%] | .0433 (.2377) [28%] |
| Sample Size | 275 | 357 | 244 |

[a]Numbers in parentheses are the apparent error rates. Because the bias is negative, our best estimate of the true error rate is the sum of the two numbers—e.g., .0344 + .3855 = .4199.

[b]Our best estimate of percentage error rate, the sum of the above numbers converted to percent, if the rule was tried on the independent sample, is in brackets.

Table 15 shows that the estimated bias in the apparent error rate is roughly 10 percent (e.g., .0344/.3855 = .09) for the two rules with lowest discriminatory power—the official record rule and the self-report official record rule. The full self-report rule has an estimated bias in the apparent error rate between 13 and 18 percent depending on the state. Thus, our estimate of the true error (misclassification) rate for our best rule (Texas) is 19 percent (.0220 + .1681) rather than 17 percent. Although the error rates observed in the RAND Inmate Survey data would be noticeably changed if the rules were applied to independent samples, the qualitative conclusions about the rules remain unchanged. We give error rates in Table 15, but they easily translate into RIOC using the formula $(p - p_c)/(1 - p_c)$ where $p$ is one minus the error rate and $p_c$ is the proportion of cases correctly labeled using a chance rule. Appendix F gives the Efron formula and a brief discussion of the rationale underlying the bias estimate.

The reason for investigating the downward bias in the apparent error rate is to guard against overfitting. That is, if the estimated bias is large, it would be prudent to use a prediction equation with fewer predictor variables. (Indeed, the larger Model 3 estimates of bias compared with Model 1 estimates reflect the fact that the more predictor variables, the more the tendency to overfit.) It is reassuring that overfitting does not appear to be a major problem for any of the discrimination rules.

# V. CONCLUSIONS

Our purpose in undertaking this research was to see how successfully one could tailor statistical discrimination methods to the problem of using official record information to distinguish between high-rate serious offenders and other offenders. In undertaking this project using self-report data from the RAND Second Inmate Survey, we were motivated by the lack of success that earlier researchers had had in addressing related questions (Chaiken and Chaiken, 1982). The discriminant rules that we developed confirm conclusions from other analysis of these data: Available official record information, particularly on arrests and convictions, contributes only marginally to identifying those inmates who have engaged in high-rate serious commission of crimes. We elaborate on the reasons for this conclusion below and also comment on its policy implications.

## STATISTICAL CONCLUSIONS

Our first step in developing an operational definition of high-rate serious offender is to estimate individual offenders' self-reported crime commission rates for the various crimes. Chaiken and Chaiken (1982) and other workers following them used the "annualized crime commission rate" defined in the obvious way: the number of reported crimes of that type divided by the amount of unincarcerated time during the measurement period. By this definition, among the survey respondents with extremely high estimated rates, some have very short periods of unincarcerated time during their measurement periods. One might argue that an offender who displayed a prodigious rate of criminal activity for a short period of time (say one to four months) could not (or would not) sustain this rate for an entire year. If so, the calculated "annualized crime commission rate" overrepresents the number of crimes he would commit if left unincarcerated for a year. On examining the data, we found evidence of "spurting" behavior; therefore we developed an adjusted estimate of each offender's annual crime commission rate that takes into account the variation throughout the year of an individual offender's crime commission rate for a particular crime. The result of this adjustment is to reduce the more extreme (above the 90th percentile) annualized rates considerably—compare our App. A with Appendix A of Chaiken and Chaiken (1982). The adjustment makes almost no difference in who is or is not classified as a high-rate serious offender according to any of the 12 definitions we used. This result is not surprising, because our high-rate serious definitions use either 70th or 80th percentiles as cutoffs; and the adjustment has its greatest influence on estimates of annualized rates above these cutoffs. We used the adjusted rates in our definitions of high-rate serious offenders.

The second step in developing operational definitions of high-rate serious offenders was exploring a variety of a priori sensible ways of specifying which crime types, committed alone or in combination, should be considered as constituting serious criminal behavior. Given a definition of serious, all offenders whose adjusted annual rate estimates are above a threshold for those crimes are deemed high-rate serious. Our six interpretations in Table 2, while not logically nested in all cases, range from classifying 3 to 15 percent of offenders in a state as high-rate serious to a high of 15 to 35 percent. We were somewhat surprised and disappointed that none of these definitions appeared to capture a natural division between the "really bad

guys" and the other offenders. Indeed, all of our definitions of "high-rate serious offender" had similar properties with respect to which explanatory variables discriminated between them. Thus, after we excluded the two definitions that yielded too few high-rate serious offenders for reliable statistical analyses (80th percentile of interpretations 1 and 2 in Table 3), our preferred definition was a somewhat arbitrary choice.

Our search for discriminant rules based on an offender's official record information that would reliably label high-rate serious offenders correctly was unsuccessful. Our best discriminant rules based on official record information did about 20 percent better than a "chance rule"[1] in correctly labeling high-rate serious offenders using our preferred definition. And the explanatory variables capturing the official record arrest and conviction information available in the RAND Survey has a modest but statistically significant relationship with being a high-rate serious offender in each of the three states. This experience is consistent with the results of Chaiken and Chaiken (1982) who did not use statistical methods specifically tailored to the discrimination problem.

When we use the self-report versions of the official record explanatory variables in an attempt to overcome the limitations of possibly low quality available official record data, the situation improves slightly. This discriminant rule correctly labels high-rate serious offenders 25 percent better than a chance rule. This finding is consistent with Chaiken and Chaiken (1982) who detected a slight improvement in estimating robbery rates when using self-report versions of official record information compared with the official record information.

Finally, as a benchmark we develop our "best" discriminant rule based on all self-report data. The improvement in correctly identifying high-rate serious offenders is substantial— almost 40 percent better than a chance rule. The improvement is due in large part to variables capturing aspects of the offender's juvenile period (crime, commitment to state facility, heroin abuse, and high school graduation) and the offender's social circumstances (employment, substance abuse) during the measurement period. Including explanatory variables on the offender's race in addition to the above variables did not improve the correct identification rate of the discriminant rule. Although the precise performance measure values of our discriminant rule vary with the definition of high-rate serious offender, the patterns are as described above.

## POLICY CONCLUSIONS

This research was, in part, motivated by the debate surrounding selective incapacitation (Blackmore and Welsh, 1984; Chaiken and Chaiken, 1985; Cohen, 1983; Fischer, 1984a; Forst, 1983; Greenwood and Abrahamse, 1982; Spelman, 1986; von Hirsch, 1984, 1985; von Hirsch and Gottfredson, 1984). Our research does not address the question of how predictive past behavior is of future commissions of crime. Prospective prediction is required for a direct test of any selective incapacitation sentencing, probation, or parole policy. Recently, Greenwood and Turner (1987) and Klein and Caggiano (1986) have carried out studies of this nature. However, if one can assume that trends in offenders' crime commission rates change only slowly over time, looking at the relationship between concurrent arrests or convictions and self-reported crime commissions is relevant to prediction of future offenses.

Because some workers have achieved respectable power with discriminant rules aimed at recidivism (e.g., Fischer, 1984a) we expected better success than we actually achieved in

---

[1]See Table 9. The RIOC criterion we use is defined as $(P - P_c)/(1 - P_c)$ where $P$ and $P_c$ are the proportion of high-rate serious offenders correctly labeled for our rule and a chance rule respectively.

discriminating between offenders who are and are not high-rate serious using official record data. We detected no evidence that the relationship in the RAND Second Inmate Survey between available official record variables, or indeed any set of explanatory variables, and being a high-rate serious offender is strong enough to be of practical use for many criminal justice policy purposes. However, the juvenile records are potentially available from official sources. The substantial improvement that our full self-report model rule made in discriminatory power gives some promise that using discriminant rules with carefully recorded juvenile record information may improve identification of high-rate serious offenders.

# Appendix A

# TABULATIONS OF ADJUSTED CRIME
# COMMISSION RATES

by Jan M. Chaiken

As described in Sec. II, the annualized crime rates previously calculated for respondents to RAND's Second Inmate Survey were adjusted in two ways for the present study: The "minimum" estimate of each individual's crime rate[1] was used rather than the average of minimum and maximum estimates, and a model was applied to adjust for the duration of the respondent's measurement period. These adjusted rates were used in determining which respondents to RAND's Second Inmate Survey were high-rate serious offenders.

The following tables give distributions and quantiles of the adjusted annualized crime commission rates for selected crime types. Tables A.1 through A.10 give statistics for California, Michigan, and Texas prison inmates, separately and combined. These tables are directly comparable to the tables for the corresponding crime types in Appendix A of Chaiken and Chaiken (1982), which give the *unadjusted* annualized rates for approximately the same subgroups of respondents.[2]

For most crime types, the quantiles (25th percentile, median, 75th percentile, and 90th percentile) of the adjusted crime commission rates are approximately 25 percent lower than in the corresponding tables of unadjusted annual crime commission rates published in 1982. The reduction tends to be somewhat larger than 25 percent in California and less than 25 percent in Michigan. These differences are caused primarily by our adoption of the minimum estimate of each individual's crime commission rate, not by our adjustment for duration of measurement period.

Although the adjustment for duration of the measurement period changes individuals' estimated crime rates substantially, it increases the estimate for some respondents and reduces it for others, thus leaving the overall distribution approximately the same. This explains why the means of the annual commission rates reported here are typically the same as, or perhaps 5–10 percent lower than, the "minimum" estimate of the mean rate reported by Chaiken and Chaiken (1982). For the crimes of business robbery and fraud, however, the adjustment for duration of the measurement period increases the mean rate over the minimum estimate published in 1982.

The label "rel" in Tables A.1 to A.10 refers to the *relative* percent of respondents whose adjusted crime rates fall in the indicated interval. The label "cum" refers to the *cumulative* percent of respondents whose adjusted crime rate is below the indicated annual rate.

---

[1]Calculation of the minimum estimate is described in Appendix B of Chaiken and Chaiken (1982).

[2]Prisoner respondents are included in the tables here only if official record information had also been collected for them; but the tables in Appendix A of Chaiken and Chaiken (1982) were calculated from data for all prisoner respondents, whether or not official record data had been collected for them. For California prisoners, our tables are based on 340 out of 357 respondents included in Appendix A of Chaiken and Chaiken (1982); for Michigan prisoners, 363 out of 422 respondents; for Texas prisoners, 583 out of 601 respondents. The combined total is 1286 prisoner respondents with official record data, out of 1380 prisoner respondents, or 93.2 percent.

Table A.1

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II: BURGLARY

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 54.2 | 45.4 | 46.8 | 44.8 |
| For actives:[b] | | | | |
| 25th percent[c] | 3.0 | 2.1 | 1.3 | 2.0 |
| Median | 6.8 | 5.6 | 3.1 | 4.8 |
| 75th percent | 84.9 | 41.9 | 7.6 | 21.8 |
| 90th percent | 327.5 | 333.6 | 82.7 | 186.7 |
| Mean | 114.6 | 107.2 | 30.4 | 77.5 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 5.3 | 5.3 | 7.6 | 7.6 | 20.7 | 20.7 | 12.3 | 12.3 |
| < 2 | 9.9 | 15.2 | 16.7 | 24.2 | 14.6 | 35.3 | 13.6 | 25.9 |
| < 3 | 9.9 | 25.1 | 12.1 | 36.4 | 13.2 | 48.5 | 11.8 | 37.8 |
| < 4 | 4.7 | 29.8 | 5.3 | 41.7 | 6.0 | 54.5 | 5.4 | 43.1 |
| < 5 | 9.9 | 39.8 | 4.5 | 46.2 | 9.4 | 63.9 | 8.4 | 51.5 |
| < 10 | 14.0 | 53.8 | 15.2 | 61.4 | 14.7 | 78.5 | 14.6 | 66.1 |
| < 20 | 12.3 | 66.1 | 6.8 | 68.2 | 5.3 | 83.8 | 8.0 | 74.0 |
| < 30 | 2.9 | 69.0 | 6.8 | 75.0 | 3.2 | 87.0 | 4.0 | 78.1 |
| < 40 | 1.2 | 70.2 | — | 75.0 | — | 87.0 | 0.4 | 78.5 |
| < 50 | — | 70.2 | 0.8 | 75.8 | 0.5 | 87.5 | 0.4 | 78.8 |
| < 100 | 7.0 | 77.2 | 4.5 | 80.3 | 2.8 | 90.3 | 4.6 | 83.5 |

[a]The percentage of respondents who reported committing burglary. This figure is taken unadjusted from App. A, Chaiken and Chaiken (1982)..

[b]The remainder of the table is based on respondents who said they committed burglary and provided data on their rates of committing burglary. The sample sizes are California, 171; Michigan, 132; and Texas, 215.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.2

## ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II: BUSINESS ROBBERY

|  | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 34.5 | 25.9 | 16.0 | 18.6 |
| For actives:[b] |  |  |  |  |
| 25th percent[c] | 1.3 | 0.7 | 0.6 | 0.8 |
| Median | 3.9 | 2.6 | 1.3 | 2.3 |
| 75th percent | 15.4 | 7.5 | 3.7 | 7.4 |
| 90th percent | 66.3 | 23.1 | 12.0 | 39.2 |
| Mean | 19.1 | 13.4 | 4.8 | 13.2 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
|  | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 15.3 | 15.3 | 32.1 | 32.1 | 41.6 | 41.6 | 28.1 | 28.1 |
| < 2 | 19.8 | 35.1 | 9.5 | 41.7 | 19.0 | 60.6 | 16.4 | 44.5 |
| < 3 | 9.0 | 44.1 | 9.5 | 51.2 | 8.6 | 69.2 | 9.1 | 53.6 |
| < 4 | 8.1 | 52.3 | 10.7 | 61.9 | 7.9 | 77.1 | 8.8 | 62.4 |
| < 5 | 5.4 | 57.7 | 3.6 | 65.5 | 2.5 | 79.6 | 4.0 | 66.4 |
| < 10 | 15.3 | 73.0 | 17.9 | 83.3 | 7.7 | 87.3 | 13.9 | 80.3 |
| < 20 | 3.6 | 76.6 | 4.8 | 88.1 | 7.6 | 94.9 | 5.1 | 85.4 |
| < 30 | 0.9 | 77.5 | 3.6 | 91.7 | 1.4 | 96.3 | 1.9 | 87.3 |
| < 40 | 4.5 | 82.0 | 1.2 | 92.9 | 1.2 | 97.5 | 2.5 | 89.8 |
| < 50 | 4.5 | 86.5 | 2.4 | 95.2 | 1.2 | 98.7 | 2.9 | 92.7 |
| < 100 | 10.8 | 97.3 | 2.4 | 97.6 | 1.3 | 100.0 | 5.5 | 98.2 |

[a]The percentage of respondents who reported committing business robbery. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed business robbery and provided data on their rates of committing business robbery. The sample sizes are California, 111; Michigan, 84; and Texas, 78.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.3

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II:
PERSON ROBBERY

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 29.6 | 26.2 | 16.9 | 20.8 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.8 | 1.7 | 0.8 | 1.4 |
| Median | 4.6 | 4.0 | 1.9 | 3.4 |
| 75th percent | 11.2 | 8.3 | 4.3 | 7.5 |
| 90th percent | 49.9 | 140.1 | 7.6 | 49.9 |
| Mean | 42.1 | 99.0 | 6.6 | 50.0 |

| | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| Distribution for actives: | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 10.5 | 10.5 | 14.8 | 14.8 | 32.6 | 32.6 | 18.5 | 18.5 |
| < 2 | 16.8 | 27.4 | 13.6 | 28.4 | 16.3 | 49.0 | 15.6 | 34.1 |
| < 3 | 6.3 | 33.7 | 7.4 | 35.8 | 12.7 | 61.6 | 8.5 | 42.6 |
| < 4 | 12.6 | 46.3 | 14.8 | 50.6 | 9.5 | 71.1 | 12.4 | 55.1 |
| < 5 | 9.5 | 55.8 | 8.6 | 59.3 | 8.3 | 79.5 | 8.9 | 63.9 |
| < 10 | 17.9 | 73.7 | 19.8 | 79.0 | 10.8 | 90.3 | 16.4 | 80.3 |
| < 20 | 7.4 | 81.1 | 1.2 | 80.2 | 6.6 | 96.9 | 5.2 | 85.5 |
| < 30 | 5.3 | 86.3 | 1.2 | 81.5 | — | 96.9 | 2.4 | 87.9 |
| < 40 | 2.1 | 88.4 | 1.2 | 82.7 | — | 96.9 | 1.2 | 89.1 |
| < 50 | 3.2 | 91.6 | 1.2 | 84.0 | — | 96.9 | 1.6 | 90.7 |
| < 100 | 2.1 | 93.7 | 1.2 | 85.2 | 1.6 | 98.4 | 1.7 | 92.3 |

[a]The percentage of respondents who reported committing robbery of persons. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed robbery of persons and provided data on their rates of committing robbery of persons. The sample sizes are California, 95; Michigan, 81; and Texas, 73.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.4

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II:
BUSINESS PLUS PERSON ROBBERY

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 48.6 | 37.6 | 25.3 | 30.1 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.8 | 1.4 | 0.7 | 1.3 |
| Median | 4.7 | 4.6 | 2.1 | 3.8 |
| 75th percent | 21.2 | 12.8 | 5.6 | 11.5 |
| 90th percent | 75.2 | 142.8 | 14.1 | 58.0 |
| Mean | 38.9 | 76.2 | 7.4 | 40.9 |

| | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| Distribution for actives: | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 11.5 | 11.5 | 20.0 | 20.0 | 31.8 | 31.8 | 20.1 | 20.1 |
| < 2 | 17.2 | 28.7 | 10.0 | 30.0 | 16.5 | 48.3 | 14.8 | 34.9 |
| < 3 | 9.6 | 38.2 | 5.8 | 35.8 | 9.6 | 57.9 | 8.4 | 43.4 |
| < 4 | 7.0 | 45.2 | 9.2 | 45.0 | 7.6 | 65.5 | 7.8 | 51.2 |
| < 5 | 6.4 | 51.6 | 5.8 | 50.8 | 7.0 | 72.6 | 6.4 | 57.6 |
| < 10 | 16.6 | 68.2 | 18.3 | 69.2 | 11.1 | 83.7 | 15.5 | 73.1 |
| < 20 | 6.4 | 74.5 | 11.7 | 80.8 | 9.1 | 92.8 | 8.8 | 81.9 |
| < 30 | 3.2 | 77.7 | 2.5 | 83.3 | 2.8 | 95.6 | 2.9 | 84.7 |
| < 40 | 5.1 | 82.8 | 1.7 | 85.0 | 0.8 | 96.4 | 2.8 | 87.6 |
| < 50 | 2.5 | 85.4 | 0.8 | 85.8 | 0.8 | 97.2 | 1.5 | 89.0 |
| < 100 | 6.4 | 91.7 | 2.5 | 88.3 | 1.8 | 99.0 | 3.8 | 92.9 |

[a]The percentage of respondents who reported committing either business robbery or robbery of persons. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed either business robbery or robbery of persons and provided data on their rates of committing either business robbery or robbery of persons. The sample sizes are California, 157; Michigan, 120; and Texas, 116.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.5

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II: ASSAULT

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 46.6 | 33.6 | 25.6 | 29.5 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.4 | 1.4 | 0.8 | 1.0 |
| Median | 3.0 | 2.6 | 1.4 | 2.1 |
| 75th percent | 9.8 | 6.0 | 3.2 | 5.7 |
| 90th percent | 18.0 | 11.9 | 6.4 | 13.6 |
| Mean | 7.5 | 5.2 | 3.3 | 5.5 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 15.4 | 15.4 | 14.5 | 14.5 | 36.5 | 36.5 | 22.0 | 22.0 |
| < 2 | 20.1 | 35.6 | 26.4 | 40.9 | 24.4 | 60.9 | 23.3 | 45.4 |
| < 3 | 12.8 | 48.3 | 12.7 | 53.6 | 9.1 | 70.0 | 11.6 | 56.9 |
| < 4 | 8.1 | 56.4 | 7.3 | 60.9 | 9.6 | 79.7 | 8.3 | 65.3 |
| < 5 | 5.4 | 61.7 | 12.7 | 73.6 | 3.2 | 82.8 | 6.8 | 72.0 |
| < 6 | 2.7 | 64.4 | 0.9 | 74.5 | 4.0 | 86.8 | 2.6 | 74.6 |
| < 7 | 6.0 | 70.5 | 4.5 | 79.1 | 3.4 | 90.2 | 4.8 | 79.4 |
| < 8 | 0.7 | 71.1 | 2.7 | 81.8 | 1.7 | 91.9 | 1.6 | 81.0 |
| < 9 | 2.7 | 73.8 | 3.6 | 85.5 | — | 91.9 | 2.1 | 83.1 |
| < 10 | 1.3 | 75.2 | 2.7 | 88.2 | — | 91.9 | 1.3 | 84.4 |
| < 20 | 16.8 | 91.9 | 8.2 | 96.4 | 5.0 | 96.8 | 10.5 | 94.8 |

[a]The percentage of respondents who reported committing assault. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed assault and provided data on their rates of committing assault. The sample sizes are California, 149; Michigan, 110; and Texas, 124.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.6

## ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II: THEFT OTHER THAN AUTO

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 41.6 | 39.7 | 36.4 | 38.0 |
| For actives:[b] | | | | |
| 25th percent[c] | 3.2 | 1.8 | 1.6 | 2.1 |
| Median | 8.5 | 4.4 | 4.7 | 5.3 |
| 75th percent | 61.8 | 53.6 | 24.3 | 51.5 |
| 90th percent | 412.6 | 317.9 | 256.3 | 291.2 |
| Mean | 159.7 | 79.3 | 153.9 | 134.1 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 8.3 | 8.3 | 16.1 | 16.1 | 16.1 | 16.1 | 13.7 | 13.7 |
| < 2 | 7.6 | 15.9 | 11.3 | 27.4 | 12.8 | 28.9 | 10.8 | 24.5 |
| < 3 | 6.8 | 22.7 | 8.1 | 35.5 | 7.7 | 36.5 | 7.5 | 32.0 |
| < 4 | 12.1 | 34.8 | 11.3 | 46.8 | 8.5 | 45.1 | 10.4 | 42.4 |
| < 5 | 2.3 | 37.1 | 8.9 | 55.6 | 5.4 | 50.5 | 5.5 | 47.9 |
| < 10 | 18.2 | 55.3 | 8.9 | 64.5 | 14.1 | 64.6 | 13.9 | 61.7 |
| < 50 | 17.4 | 72.7 | 10.5 | 75.0 | 11.0 | 75.7 | 12.8 | 74.6 |
| < 100 | 5.3 | 78.0 | 5.6 | 80.6 | 5.6 | 81.3 | 5.5 | 80.1 |
| < 200 | 4.5 | 82.6 | 8.1 | 88.7 | 5.4 | 86.6 | 5.9 | 86.0 |
| < 300 | 3.8 | 86.4 | 0.8 | 89.5 | 3.5 | 90.1 | 2.8 | 88.8 |
| < 500 | 5.3 | 91.7 | 6.5 | 96.0 | 3.6 | 93.7 | 4.9 | 93.7 |

[a]The percentage of respondents who reported committing theft other than auto. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed theft other than auto and provided data on their rates of committing theft other than auto. The sample sizes are California, 132; Michigan, 124; and Texas, 169.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.7

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II:
AUTO THEFT

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 24.3 | 23.2 | 18.8 | 20.4 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.1 | 1.9 | 0.8 | 1.0 |
| Median | 2.6 | 3.7 | 1.5 | 2.2 |
| 75th percent | 7.1 | 48.0 | 3.6 | 6.6 |
| 90th percent | 69.6 | 379.4 | 14.9 | 76.3 |
| Mean | 28.7 | 213.5 | 8.2 | 74.1 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 22.2 | 22.2 | 13.3 | 13.3 | 28.8 | 28.8 | 22.1 | 22.1 |
| < 2 | 18.1 | 40.3 | 11.7 | 25.0 | 32.5 | 61.3 | 21.6 | 43.7 |
| < 3 | 12.5 | 52.8 | 15.0 | 40.0 | 6.3 | 67.7 | 10.9 | 54.6 |
| < 4 | 4.2 | 56.9 | 16.7 | 56.7 | 9.3 | 77.0 | 9.7 | 64.3 |
| < 5 | 5.6 | 62.5 | 3.3 | 60.0 | — | — | 2.9 | 67.1 |
| < 10 | 18.1 | 80.6 | 10.0 | 70.0 | 9.3 | 86.3 | 12.5 | 79.7 |
| < 20 | 2.8 | 83.3 | 3.3 | 73.3 | 8.1 | 94.5 | 4.9 | 84.6 |
| < 30 | 4.2 | 87.5 | — | — | — | — | 1.4 | 86.0 |
| < 40 | — | — | 1.7 | 75.0 | — | — | 0.5 | 86.5 |
| < 50 | 1.4 | 88.9 | — | — | 1.5 | 96.0 | 1.0 | 87.5 |
| < 100 | 4.2 | 93.1 | 6.7 | 81.7 | 1.2 | 97.2 | 3.8 | 91.3 |

[a]The percentage of respondents who reported committing auto theft. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed auto theft and provided data on their rates of committing auto theft. The sample sizes are California, 72; Michigan, 60; and Texas, 75.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.8

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II:
FORGERY PLUS CREDIT CARDS

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 28.4 | 14.1 | 21.5 | 20.9 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.4 | 1.7 | 1.5 | 1.5 |
| Median | 3.5 | 4.7 | 3.3 | 3.4 |
| 75th percent | 15.0 | 51.6 | 9.7 | 11.7 |
| 90th percent | 89.6 | 524.0 | 66.4 | 84.0 |
| Mean | 52.5 | 105.6 | 28.7 | 51.5 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 17.2 | 17.2 | 16.7 | 16.7 | 18.9 | 18.9 | 17.9 | 17.9 |
| < 2 | 13.8 | 31.0 | 11.9 | 28.6 | 12.5 | 31.3 | 12.9 | 30.7 |
| < 3 | 12.6 | 43.7 | 11.9 | 40.5 | 13.3 | 44.6 | 12.8 | 43.5 |
| < 4 | 11.5 | 55.2 | 9.5 | 50.0 | 15.6 | 60.2 | 13.0 | 56.5 |
| < 5 | 10.3 | 65.5 | — | 50.0 | 5.0 | 65.2 | 6.1 | 62.6 |
| < 10 | 9.2 | 74.7 | 11.9 | 61.9 | 9.4 | 74.6 | 9.8 | 72.3 |
| < 20 | 4.6 | 79.3 | 7.1 | 69.0 | 8.5 | 83.1 | 6.8 | 79.1 |
| < 30 | 1.1 | 80.5 | — | 69.0 | — | 83.1 | 0.4 | 79.6 |
| < 40 | 2.3 | 82.8 | 4.8 | 73.8 | 0.9 | 83.9 | 2.1 | 81.7 |
| < 50 | — | 82.8 | — | 73.8 | 1.8 | 85.8 | 0.8 | 82.5 |
| < 100 | 10.3 | 93.1 | 11.9 | 85.7 | 6.6 | 92.3 | 8.9 | 91.4 |

[a]The percentage of respondents who reported committing forgery or credit card swindles. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed forgery or credit card swindles and provided data on their rates of committing forgery or credit card swindles. The sample sizes are California, 87; Michigan, 42; and Texas, 103.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.9

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II: FRAUD

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 19.3 | 16.1 | 14.2 | 15.2 |
| For actives:[b] | | | | |
| 25th percent[c] | 1.0 | 0.9 | 0.9 | 1.0 |
| Median | 3.2 | 2.4 | 2.9 | 2.9 |
| 75th percent | 25.2 | 7.7 | 6.9 | 8.7 |
| 90th percent | 89.9 | 66.6 | 34.2 | 66.6 |
| Mean | 48.1 | 20.8 | 18.8 | 28.5 |

| | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| Distribution for actives: | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 1 | 21.1 | 21.1 | 27.5 | 27.5 | 28.2 | 28.2 | 25.8 | 25.8 |
| < 2 | 15.8 | 36.8 | 19.6 | 47.1 | 9.7 | 38.0 | 14.4 | 40.2 |
| < 3 | 10.5 | 47.4 | 13.7 | 60.8 | 12.6 | 50.5 | 12.3 | 52.4 |
| < 4 | 5.3 | 52.6 | 7.8 | 68.6 | 13.3 | 63.8 | 9.3 | 61.7 |
| < 5 | 5.3 | 57.9 | — | 68.6 | 1.3 | 65.1 | 2.2 | 63.8 |
| < 10 | 8.8 | 66.7 | 9.8 | 78.4 | 13.3 | 78.4 | 10.9 | 74.7 |
| < 20 | 8.8 | 75.4 | 3.9 | 82.4 | 7.9 | 86.3 | 7.1 | 81.8 |
| < 30 | — | 75.4 | — | 82.4 | — | 86.3 | — | 81.8 |
| < 40 | 5.3 | 80.7 | 2.0 | 84.3 | 4.3 | 90.6 | 3.9 | 85.8 |
| < 50 | 5.3 | 86.0 | 2.0 | 86.3 | — | 90.6 | 2.2 | 88.0 |
| < 100 | 5.3 | 91.2 | 11.8 | 98.0 | 2.6 | 93.2 | 6.0 | 94.0 |

[a]The percentage of respondents who reported committing fraud. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who said they committed fraud and provided data on their rates of committing fraud. The sample sizes are California, 57; Michigan, 51; and Texas, 74.

[c]Percentiles shown at the top of the table are defined by interpolating between actual data points; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

Table A.10

ADJUSTED CRIME COMMISSION RATES—INMATE SURVEY II:
DRUG DEALING

| | California Prison | Michigan Prison | Texas Weighted Prison | Total |
|---|---|---|---|---|
| Percent active[a] | 54.5 | 41.4 | 34.6 | 41.4 |
| For actives:[b] | | | | |
| 25th percent[c] | 7.4 | 5.7 | 5.2 | 6.2 |
| Median | 78.6 | 137.5 | 29.4 | 67.1 |
| 75th percent | 547.0 | 451.5 | 336.3 | 403.4 |
| 90th percent | 2669.7 | 2636.5 | 2019.1 | 2487.3 |
| Mean | 849.9 | 1011.0 | 655.7 | 826.9 |

| Distribution for actives: | Percent | | Percent | | Percent | | Percent | |
|---|---|---|---|---|---|---|---|---|
| | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. | Rel. | Cum. |
| < 5 | 18.5 | 18.5 | 23.3 | 23.3 | 23.7 | 23.7 | 21.7 | 21.7 |
| < 10 | 11.0 | 29.5 | 11.3 | 34.6 | 16.0 | 39.7 | 12.8 | 34.5 |
| < 50 | 17.3 | 46.8 | 12.0 | 46.6 | 13.3 | 52.9 | 14.4 | 48.9 |
| < 100 | 4.6 | 51.4 | 0.8 | 47.4 | 4.8 | 57.8 | 3.6 | 52.5 |
| < 500 | 22.5 | 74.0 | 27.8 | 75.2 | 21.1 | 78.9 | 23.5 | 76.0 |
| < 1000 | 7.5 | 81.5 | 5.3 | 80.5 | 6.2 | 85.1 | 6.4 | 82.5 |
| < 3000 | 9.8 | 91.3 | 12.8 | 93.2 | 10.4 | 95.4 | 10.8 | 93.3 |

[a]The percentage of respondents who reported dealing drugs. This figure is taken unadjusted from Appendix A, Chaiken and Chaiken (1982).

[b]The remainder of the table is based on respondents who reported dealing drugs and provided data on their rates of dealing drugs. The sample sizes are California, 173; Michigan, 133; and Texas, 163.

[c]Percentiles shown at the top of the table are defined by interpolation between actual data prints; e.g., if data occur in a 10 percent chunk corresponding to the 85th and 95th percentile, the 90th percentile is estimated as half way between two adjacent data points. This definition accounts for possible slight inconsistencies between the percentiles and the distribution at the bottom of the table.

# Appendix B

# DISPARITIES IN CRIME RATE ESTIMATES
# FROM TWO RAND INMATE SURVEYS

by Jan M. Chaiken

RAND's First Inmate Survey, in 1976, was an anonymous written survey of California prisoners.[1] The Second Inmate Survey, carried out in late 1978 and early 1979, covered three states (Michigan and Texas as well as California) and included inmates of county jails as well as state prisons.[2]

The questionnaire instrument used in the second survey was a refinement and elaboration of the first instrument, but both included questions about the numbers of crimes respondents had committed in a period preceding their incarcerations. The answers to these questions were converted into estimates of each respondent's annual crime commission rates—the number of crimes (of each of several types) that the respondent reported per year of unincarcerated time.[3] The published distributions differed remarkably between the two surveys. For example, the comparisons in Table B.1 show estimated mean rates from the second survey ranging between 1.6 and approximately 20 times as high as estimated rates for the same type of crime from the first survey. The estimated prevalences (percent of the cohort that committed the crime) are much closer but also show some substantial differences between the surveys.

Many possible explanations of the disparities can be proposed, including the following:

1. *Sample composition.* The sampling design differed between the two surveys. In 1976 the respondents presumably constituted approximately a random sample of incarcerated prisoners. In 1979, the sample was chosen to be representative of an incoming cohort,[4] and indeed the respondents' characteristics are distributed like those of an incoming cohort (Peterson et al., 1982, p. 59). A simple random sample of incarcerated prisoners would have a greater proportion of people serving long sentences (e.g., murderers or kidnappers) than the proportion in an incoming cohort. To adjust for the sampling method, a mathematical model was used to estimate weights for individual respondents in the 1976 survey in order to estimate what crime commission rates would be for an incoming cohort. The 1979 data were processed without using weights (i.e., each respondent counted the same as any other respondent in computing crime commission estimates).

2. *Cohort and selection effects.* By dint of the passage of three years, the criminal behavior of individuals in prison in 1979 could in fact have been different from that of 1976 prison inmates.

---

[1]See Peterson and Braiker (1981).

[2]For a description of the survey, see Peterson et al. (1982). The California prisoners were surveyed in 1979.

[3]The respondent's unincarcerated time during the survey reference period could have ranged from one month to several years.

[4]An incoming cohort is a group of people who begin serving a prison sentence during a given period of time, such as a year or a month. The 1979 sample was not a subset of an incoming cohort but rather simulated an incoming cohort.

### Table B.1

#### ESTIMATED CRIMINAL BEHAVIOR FOR AN INCOMING COHORT OF CALIFORNIA PRISONERS

| Crime Type | Percent of Cohort Committing the Crime | | Mean Annual Rate for Cohort Members Who Commit the Crime | |
|---|---|---|---|---|
| | 1976 Survey | 1979 Survey | 1976 Survey | 1979 Survey |
| Armed robbery | 37 | 4.6 | | |
| Robbery | | 49 | | 49-74 |
| Burglary | 58 | 54 | 15.3 | 116-204 |
| Assault | 59 | 47 | 4.5 | 7.1-7.6 |
| Drug deals | 48 | 55 | 155 | 927-1681 |
| Auto theft | 32 | 24 | 5.3 | 38-102 |
| Forgery | 40 | 28 | 5.6 | 62-94 |
| Cons | 63 | | 9.5 | |
| Fraud | | 19 | | 156-202 |

SOURCES: Peterson and Braiker (1981), Table 12. Chaiken and Chaiken (1982), Appendix A.

NOTE: The crime commission data from the second inmate survey were processed to yield a minimum and maximum estimate of each crime rate for each respondent (Chaiken and Chaiken, 1982, Appendix A). The minimum usually equaled the maximum but could differ because of incomplete or ambiguous responses etc. The last column gives averages of the minimums and averages of the maximums.

3. *Definitions of crime types.* The descriptions of types of criminal behavior differed between the two surveys (Table B.2).[5] For questions on some crime types the differences are so major that we have given them distinct labels in Tables B.1 and B.2 (e.g., armed robbery vs. robbery). For others (e.g., burglary), the wording is almost identical in both surveys, but the differences between wordings could account for substantial disparities in answers for some respondents.[6]

4. *Format of survey questions.* The crime commission rates reported for the Second Inmate Survey were calculated from questions that asked respondents who had committed fewer than 11 crimes (of a given type) to report the number committed (1, 2, . . ., 10) and asked respondents with higher counts to tell their *frequencies* of commission (e.g., crimes per month, per week, or per day). By contrast, the rates reported from the first survey were based on *categorized* responses to a question about total *counts* of crimes committed during the three-year measurement period. The categories in the first survey instrument were as follows:[7]

---

[5]The wording in the Second Inmate Survey was intended to conform more closely to legal definitions, so that self-reported arrests could be compared with officially recorded arrests.

[6]For example, a person might break into cars many times (included in the 1978–79 wording for burglary), but never into houses or businesses (1976 wording).

[7]The categories differed for drug sales: 0, less than 10, less than 50, less than 100, more than 100.

## Table B.2

### COMPARISON OF WORDING OF SURVEY QUESTIONS

| 1976: First Inmate Survey | 1978–79: Second Inmate Survey |
|---|---|
| *Armed robbery*<br>Threatened someone with a weapon in order to get money or something else. | *Robbery*<br>Hold up a store, gas station, bank, taxi, or other business. Rob any persons, do any muggings, street robberies, purse snatches, or holdups in someone's house or car. |
| *Burglary*<br>Broke into a home or business in order to take something. | *Burglary*<br>Broke into a home or a car or a business in order to take something. |
| *Assault*<br>Beat or physically hurt someone badly. Cut someone with a knife or shot someone with a gun. Threatened to hurt someone with a gun, knife or other weapon. Tried to kill someone. | *Assault*<br>Assault someone, threaten someone with a weapon, shoot at someone, try to cut someone, or beat or strangle someone. |
| *Drug sales*<br>Sold hard drugs. | *Drug deals*<br>Deal in drugs: Make, sell, smuggle or move drugs. |
| *Auto theft*<br>Stole a car. | *Auto theft*<br>Steal any cars, trucks or motorcycles. |
| *Forgery*<br>Forged a check or other paper. | *Forgery*<br>Forge something, use a stolen or bad credit card, or pass a bad check. |
| *Cons*<br>Hustled or conned someone. | *Fraud*<br>Do any frauds or swindles (illegal cons) on a person, business, or the government. |

- 0
- 1–2
- 3–5
- 6–10
- More than 10, how many?

The first survey also included questions asking for frequencies.[8] The mean rate for burglary calculated from the "frequency" question in the 1976 survey was reported[9] to be 16 times the 1976 survey categorical question rate shown in Table B.1, or about 240 burglaries per year, which is even higher than the estimated burglary rate calculated from the 1978–79 survey (also shown in Table B.1).

5. *Differing measurements periods.* In the 1978–79 survey the duration of the measurement period for reporting crimes varied between 13 and 24 months. Every respondent's measurement period began in a January and ended with his arrest for the crime he was serving time for. In the 1976 survey, the measurement period was three years long (36 months) and

---

[8]The data from the frequency questions were not used in analyses of the first survey, as explained by Peterson and Braiker (1980), pp. 25–26.

[9]Peterson and Braiker (1981), note 19, p. 26.

ended when he *began serving* his current term. To the extent there is "spurting" of crime commissions or crime sprees, different length measurement periods would affect estimated rates. Also, criminal behavior between arrest and incarceration could be different from behavior just before arrest.

6. *Handling of uncertain responses.* Some respondents to the 1976 survey failed to check any of the crime count categories listed in the question; the analysts assumed these respondents intended to report zero commissions but had not understood the numeral 0. Further, the data processing program eliminated certain outlier (high) responses[10] and imputed specific values for responses "more than 10, unspecified how much."[11] For purposes of calculating mean rates, responses in categories were set equal to their midpoints (e.g., the range 6–10 was set equal to 8).

This appendix discusses our analyses related to the first four of these issues. Our results concerning the effects of different measurement periods, and researchers' handling of uncertain responses, are in Sec. II of the text.

## COMPARING FORMATS IN 1976 AND 1978–79

To help sort out the influence of those various factors, the Second Inmate Survey instrument replicated selected 1976 survey questions for categorized crime counts. These questions had the same wording and the same response categories as the 1976 survey, except that the category "more than 10" did not continue on to "how many?" (Appendix C gives the 1976 format and Appendix D the 1978–79 replication format.) The 1976 question concerning armed robbery was not repeated in the second survey in order to avoid possibly confusing the respondents about the definition of "robbery" as used in questions that would have followed the armed robbery question.

Responses to these replication questions have not been used in any earlier published analyses of the Second Inmate Survey (except for measuring internal reliability of respondents' answers). All previously reported crime commission rates for respondents to the second survey were derived from more detailed questions that were separated from the replication questions by as much as 25 pages. (See Appendix E for an example of the format of questions used to estimate respondents' crime rates.)

By comparing the two sets of responses to the 1978–79 survey we were able to assess the role of analytical decisions and instrumentation differences in creating disparities between the 1976 results and the 1979 results for California prisoners. These results, described below, show that most data processing and sample composition issues are unimportant, as are the precise details of crime definitions and the distinctions between *counts* of crimes (1, 2, . . . , 10) and *categorized counts*. We are able to confirm that a question format asking for *frequencies* yields a higher estimate than one asking for *counts*, but it is not clear which is closer to the truth.

Our analysis reported here was not limited to California prisoners but included all 2190 respondents whose crime commission rates were reported for the Second Inmate Survey. Four crime types were sufficiently similar to allow fairly direct comparisons between the two questionnaire formats: burglary, auto theft, forgery, and drug deals. For these four, the 1978–79 wording (second column of Table B.2) logically allowed for more crimes to be reported than did the replication wording (first column). Thus, it is reasonable to expect small disparities in

---

[10]Peterson, Polich, and Chaiken (1980), p. 33.

[11]Peterson, Polich, and Chaiken (1980), p. 20.

responses solely due to wording. We also compared the responses for "cons" with those for "fraud." In this case the replication wording (from 1976) was much more liberal than the new 1978–79 wording, leading one to expect higher responses to the replication question solely due to the difference in wording.

## EFFECT OF IMPUTED ZEROS FOR MISSING RESPONSES IN 1976 SURVEY

The original editing decision to impute "zero" for the missing responses to the 1976 survey was found not to play a role in the disparities between the results of the two surveys. The "new" 1978–79 format included a specific "Yes-No" question to determine whether the respondent committed the crime at all during the measurement period. (See Fig. B.1 for an example.) Some respondents left this question blank (or, in very rare instances, gave an ambiguous response—such as both "No" and "Yes"). The responses other than "Yes" or "No" constituted from 1.5 percent of the sample (for burglary and for auto theft) to 3.2 percent (for forgery). These figures are only slightly lower than the percentage of blank responses in the 1976 survey (2.4 to 3.6 percent). Moreover, the 1978–79 format of the replication question—which allowed the respondent to check a box for "more than 10" instead of requiring (as in 1976) a numerical answer to "more than 10, how many?"—did not elicit fewer blank responses. The blank responses for the replication questions in 1978–79 constituted from 5.3 percent of the sample (for burglary) to 5.9 percent (for auto theft).[12]

Table B.3 compares the answers to the new Yes-No questions for respondents who left the replication question blank with corresponding answers to the new questions for remaining respondents. On the whole, the respondents who left the replication questions blank were less likely to respond "yes" to the new questions than were other respondents. Moreover, for all five crime types studied, over 70 percent of those who left the replication question blank answered "No" to the corresponding direct question.

We analyzed the crime rates for the small number of respondents who answered "yes" to the new 1978–79 question while leaving the corresponding replication question blank. They had disproportionately low (not high) estimated crime rates.

We conclude that the number of respondents who left the 1976-style replication question blank was so small that imputing their responses as zero (rather than missing) would not substantially affect any estimates derived from these questions for the entire sample. Further, imputation of zero was more likely to be correct than not and therefore was a reasonable analytic choice.[13] Our analysis of the 1978–79 survey data does not suggest that the "true" behavior of offenders who failed to answer the 1976 survey questions about crime commissions could in any way account for disparities between the 1976 and 1979 results for California prisoners.

### Comparison of Survey Formats: 1978–79 Survey

The replication questions on the 1978–79 survey asked the respondent to check a box indicating a categorized count of the number of crimes committed during the measurement period. The new format 1978–79 questions also asked for the count of crimes, if 10 or less;

---

[12]4.0 percent of respondents to the 1978–79 survey left all the replication questions blank. For California prisoner respondents to this survey, from 3.6 to 5.1 percent of replication questions were left blank.

[13]However, in the remainder of the work reported here, we have handled them as missing.

Replication Question

...how many times did you do each of the following:

    d.  Burglary--broke into a home or business in order to take
        something

        0 ☐    1-2 ☐    3-5 ☐    6-10 ☐    More than 10 ☐

New Format (from page 16 of survey booklet)

1.  During the STREET MONTHS ON THE CALENDAR did you do any burglaries?
    (Count any time that you broke into a house or a car or a business
    in order to take something.)

        YES ☐₁         NO ☐₂ → go on to page 18

2.  In all, how many burglaries did you do?

    ☐ 11 OR MORE        ☐ 1 TO 10
       ↓             How many?

3.  Look at the total street           /_____/ Burglaries
    months on the calendar.
    During how many of those      go on to next page ⟶
    months did you do one or
    more burglaries?

    _____ Months

4.  In the months when you did burglaries,
    how often did you usually do them?

        (CHECK ONE BOX)

    EVERYDAY OR         How many    /_____/  How many days  /_____/
    ALMOST EVERYDAY ☐→per day?         → a week usually?

    SEVERAL TIMES       How many    /_____/
    A WEEK ☐→per week?

    EVERY WEEK OR      How many    /_____/
    ALMOST EVERY WEEK ☐→per month?

    LESS THAN          How many    /_____/
    EVERY WEEK ☐→per month?

Fig. B.1—Example of replication question and new format
in 1978–79 survey

## Table B.3

### DISTRIBUTION OF ANSWERS TO YES-NO CRIME COMMISSION QUESTIONS

(Percent of respondents to 1978–79 Inmate Survey)

| Crime Type | Replication Question | Answer to "Did you commit?" (new format question) | | | | Number of Respondents |
| | | No | Yes | Blank | Total | |
|---|---|---|---|---|---|---|
| Burglary | | | | | | |
| | Blank | 73 | 18 | 9 | 100 | 116 |
| | Not Blank | 59 | 40 | 1 | 100 | 2065 |
| Auto Theft | | | | | | |
| | Blank | 80 | 8 | 12 | 100 | 128 |
| | Not Blank | 82 | 17 | 1 | 100 | 2052 |
| Forgery | | | | | | |
| | Blank | 73 | 17 | 10 | 100 | 124 |
| | Not Blank | 79 | 19 | 2 | 100 | 2055 |
| Fraud/cons | | | | | | |
| | Blank | 82 | 12 | 5 | 100 | 112 |
| | Not Blank | 83 | 15 | 2 | 100 | 2068 |
| Drug deals | | | | | | |
| | Blank | 71 | 17 | 12 | 100 | 119 |

NOTE: Total sample size = 2190. Missing cases gave ambiguous or unclear responses.

however, in this format the respondent gives the count directly by writing it, not by checking a box, as shown in Fig. B.1.

For four of the five crime types being studied here, the top category of the replication questions is also "more than 10," so respondents were asked about the total count of crimes committed, either categorized or not, in the two formats. But for the fifth crime type, dealing drugs, the top category of the replication question was "more than 100." If their answer was between 11 and 100 drug deals, the respondents were asked to provide a *categorized count* of crimes in one format and a crime *frequency* in the other.

Thus a comparison of the count of crimes reported in two formats yields the following information:

- for the crimes of burglary, auto theft, and forgery: a comparison of categorized and uncategorized formats
- for the crime of fraud/cons: a determination of whether the large definitional difference of the crime types in the two formats yields observably different results
- for the crime of drug dealing: a comparison of responses using a *frequency* format vs. a *categorized count* format.

The analysis was carried out by estimating the count of crimes from the "new" format (either directly from the response or by multiplying the frequency times the appropriate time

period shown on the questionnaire)[14] and categorizing the count to coincide with the categories used in the replication questions.[15]

The resulting cumulative distributions for four of the crime types are shown in Fig. B.2. As an example, looking at the unshaded bars for burglary: 57 percent of respondents to the replication question answered "zero"; 16 percent answered "1–2," for a cumulative total of 73+ percent; 9+ percent answered "3–5," for a cumulative total of 83 percent; and 5 percent answered "6–10" for a cumulative 88 percent. The remaining 12 percent (visible as the space between the bar on the right and the top of the graph) responded "more than 10." Visually, the distributions from the two questionnaire formats appear remarkably similar for the crimes burglary, auto theft, and forgery.[16] In particular, the estimated prevalence (fraction of respondents who committed the crime at all during the measurement period) is essentially the same using either format. For the crime of drug dealing, the new format frequency question yields a similar distribution up to 50 drug deals but somewhat fewer respondents in the range 51–100 and correspondingly more in the category "over 100." The prevalence estimate is again essentially the same from either format.

By contrast, Fig. B.3 shows that the same type of comparison between two crimes, cons and frauds, that are worded substantially differently generates two distributions that are not visually close.

The general similarity of the distributions in Fig. B.2 obscures the fact that the reliability of individual responses was not high. As shown in Table B.4, the percentage of respondents whose categorized "new" format response agreed with his answer to the replication question varied between 49 and 70 percent. Only when adjacent categories are considered (e.g., "3–5" is adjacent to "1–2") does the correspondence reach high levels.

The poor reliability of the drug dealing question, indicated in Table B.4, led us to examine whether the reported count of crimes committed—which was intended to refer to the entire measurement period—could possibly have coincided with the *daily* or *monthly* rate of crimes reported on a previous page in the survey booklet. This was found to be a likely explanation in some cases. For example, of 27 people who checked "less than 10" on the replication question but had a count over 100 according to the "new" format 1978–79 question, the rate per day or per month was between 1 and 10 for 16 (59 percent) of them. For the 32 who responded "less than 50" (but not "less than 10") and whose new-format count was over 100, the daily or monthly rate was between 11 and 50 drug deals for 21 (66 percent) of them. Thus the correspondence between reported total count in one format and reported monthly or daily count in the other format is substantially better than one would expect at random.

In sum, the analysis suggests that either of the following formats will yield very similar self-reports of the number of crimes committed during a specified measurement period:
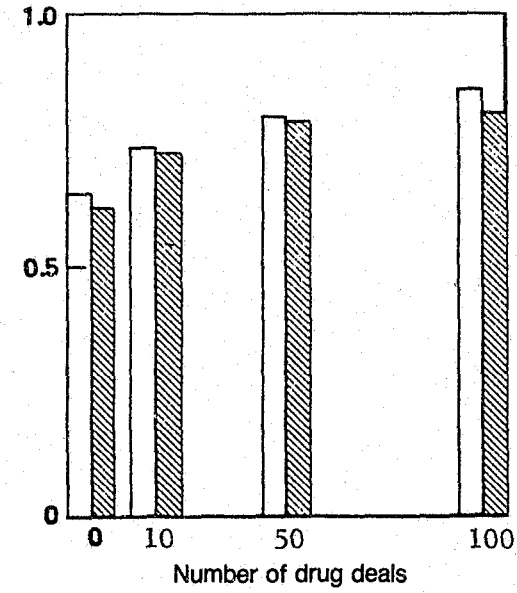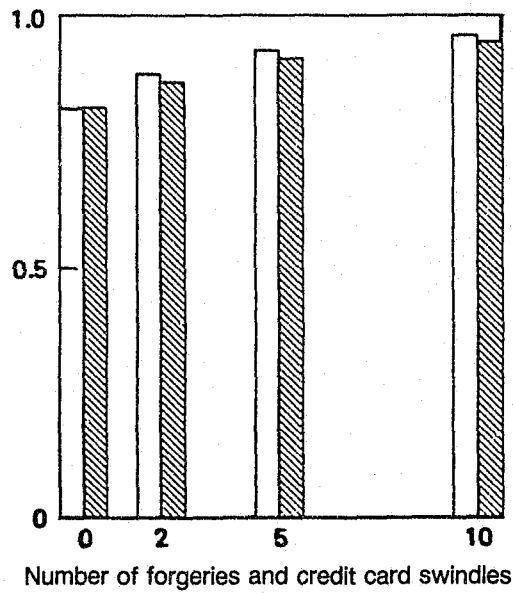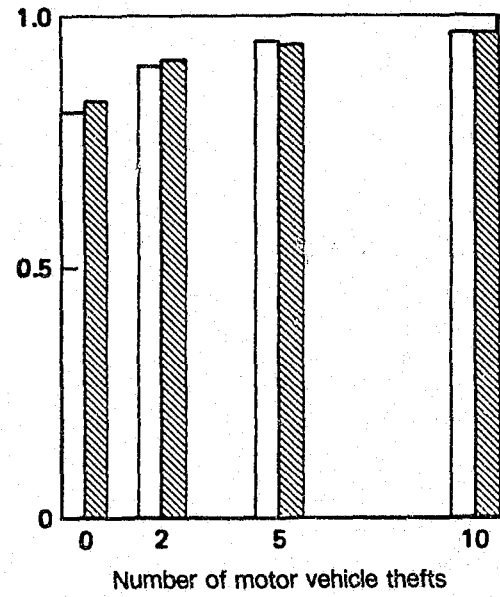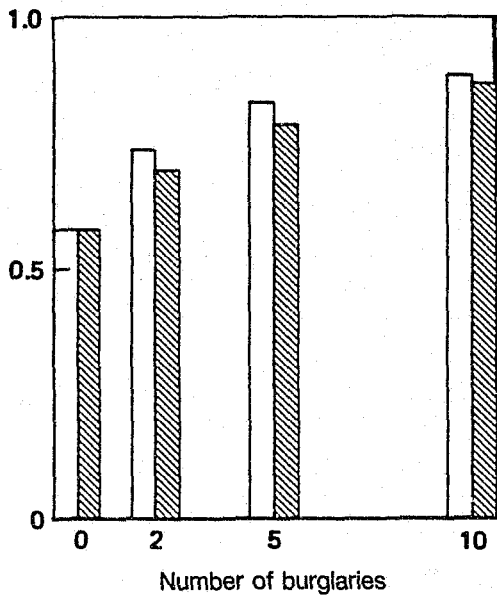
- Boxes to be checked for the categories 0, 1–2, 3–5, 6–10, and over 10.
- A Yes-No question, which if answered "No" indicates "zero"; a box to enter a number between 1 and 10; and another box to check indicating "11 or more."

---

[14]See question 3 in Fig. B.1 for an example of the time-period question.

[15]In case of an ambiguous or "range" response in the new format, we gave the respondent the benefit of the doubt. That is, we chose the permissible category closest to the response on the replication question. For each crime type, fewer than 26 responses (1.2 percent of the sample) were ambiguous in relation to the categories.

[16]Because of the large sample size, however, a $\chi^2$ test shows that the data are not consistent with the hypothesis that the pairs of empirical distributions are drawn from the same underlying distribution. The $\chi^2$ values with 4 degrees of freedom are: burglary 79, auto theft 35, forgery 24, drug dealing 64.

Fig. B.2—Cumulative distribution of the number of various crimes
committed during the measurement period

"Cons" from replication question
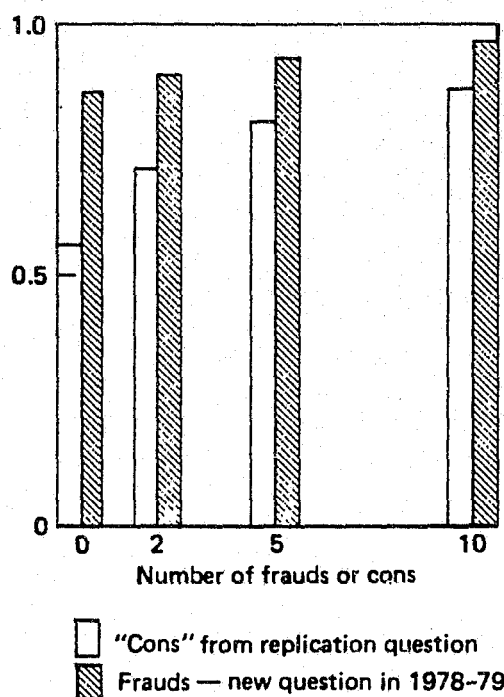
Frauds — new question in 1978–79

Fig. B.3—Cumulative distribution of the number of frauds or cons
committed during the measurement period

Table B.4

PERCENT OF RESPONDENTS WITH CLOSE CATEGORIES IN TWO FORMATS

| Crime Type | Number of Respondents with Crime Count > 0 | Percent of Respondents | |
|---|---|---|---|
| | | Same Category | Same or Adjacent Category |
| Burglary | 808 | 62 | 88 |
| Auto theft | 338 | 71 | 95 |
| Forgery | 366 | 60 | 87 |
| Drug deals | 762 | 49 | 77 |

If the reported count of crimes committed is over 50, our analysis (of the drug sales question) suggests that a format asking for a frequency (number of crimes per month, per week, or per day) will yield higher frequencies than a format that asks directly for the total count. Which of these answers is closer to the truth is unclear, because we found some indications that the reported "total count" might actually be the respondent's estimate of his daily or monthly count of crimes committed.[17]

---

[17]This conclusion is weakened because of differences in wording of the two particular questions being compared: The categorized question asked about sales of hard drugs; and the new format question asked about all kinds of drug deals, which could in reality be more numerous.

# Appendix C

## CRIME COMMISSION QUESTIONS FROM 1976 RAND INMATE SURVEY

19. During those 3 years before you started your present term, all together how many times did you do each of following:

a. Armed robbery--threatened someone with a weapon in order to get money or something else.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  58-60/

b. Totally lost your temper.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  61-62/

c. Beat or physically hurt someone badly.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many ____  63-64/

d. Hustled or conned someone.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  65-66/

e. Cut someone with a knife or shot someone with a gun.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  67-68/

f. Burglary--broke into a home or business in order to take something.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  69-71/

g. Got into a fist fight.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  72-73/

h. Forced someone to have sex with you.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  74-75/

i. Got drunk and hurt someone.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  76-77/

CARD 03

j. Threatened to hurt someone with a gun, knife or other weapon.

   0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10, how many? ____  10-11/

Go on to next page ─────────➤

k.  Tried to kill someone.

   0 ☐     1-2 ☐     3-5 ☐     6-10 ☐     More than 10, how many? ____ *12-13/*

l.  Forged a check or other paper.

   0 ☐     1-2 ☐     3-5 ☐     6-10 ☐     More than 10, how many? ____ *14-16/*

m.  Stole a car.

   0 ☐     1-2 ☐     3-5 ☐     6-10 ☐     More than 10, how many? ____ *17-19/*

n.  Sold hard drugs.

   0 ☐     Less than 10 ☐     Less than 50 ☐     Less than 100 ☐     *20-24/*

   More than 100, how many? ____

# Appendix D

# REPLICATION OF SELECTED 1976 QUESTIONS
# IN 1978–79 RAND INMATE SURVEY

14. Again look at the calendar. During the STREET MONTHS ON THE CALENDAR, altogether how many times did you do each of the following:

a. Beat or physically hurt someone badly.

  0 ☐₀   1-2 ☐₁   3-5 ☐₂   6-10 ☐₃   More than 10 ☐₄

b. Hustled or conned someone.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

c. Cut someone with a knife or shot someone with a gun.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

d. Burglary--broke into a home or business in order to take something.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

e. Threatened to hurt someone with a gun, knife or other weapon.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

f. Tried to kill someone.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

g. Forged a check or other paper.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

h. Stole a car.

  0 ☐   1-2 ☐   3-5 ☐   6-10 ☐   More than 10 ☐

i. Sold hard drugs.

  0 ☐   Less than 10 ☐   Less than 50 ☐   Less than 100 ☐   More than 100 ☐

## Appendix E

## EXAMPLE NEW FORMAT CRIME COMMISSION QUESTION
## IN 1978–79 RAND INMATE SURVEY

The next questions are also only about the STREET MONTHS ON THE CALENDAR.
Look at the calendar to help you remember what you were doing during these
months. These are months that do not have X's or lines in them.

1. During the STREET MONTHS ON THE CALENDAR did you do any burglaries?
   (Count any time that you broke into a house or a car or a business
   in order to take something.)

   YES ☐₁                    NO ☐₂ ➡ go on to page 18

2. In all, how many burglaries did you do?

   ☐ 11 OR MORE                    ☐ 1 TO 10
                                        How many?

3. Look at the total street
   months on the calendar.                      ⬡ Burglaries
   During how many of those
   months did you do one or                go on to next page ⟶
   more burglaries?

   _____ Months

4. In the months when you did burglaries,
   how often did you usually do them?

   *(CHECK ONE BOX)*
                        ⬇
   EVERYDAY OR          ☐ How many          ⬡   How many days
   ALMOST EVERYDAY        ⟶per *day?*            ⟵ a week usually? ⬡

   SEVERAL TIMES        ☐ How many          ⬡
   A WEEK                 ⟶per *week?*

   EVERY WEEK OR        ☐ How many          ⬡
   ALMOST EVERY WEEK      ⟶per *month?*

   LESS THAN            ☐ How many          ⬡
   EVERY WEEK             ⟶per *month?*

# Appendix F

# METHODOLOGY FOR ESTIMATING BIAS
# IN ERROR RATES

Efron (1986) develops a simple estimate for the downward bias of the apparent error rate for discrimination rules based on logistic regression. Efron's method is used to compute the bias estimates given in Table 15. We describe Efron's method in this appendix.

Let $p(x_i)$ be the maximum likelihood fitted logistic probability that offender i with vector of characteristics $x_i$ is high-rate serious. Our discrimination rule is defined as

$$D(i) = 1 \text{ if } p(x_i) > c ,$$

$$= 0 \text{ if } p(x_i) \leq c . \tag{F.1}$$

If $y_i = 1$ offender i is high-rate serious and $y_i = 0$ if offender i is not, then c is chosen so that the proportion of $y_i = 1$ is the same as the proportion for D(i). The apparent error rate (AER) of rule (F.1) is

$$AER = (\text{number of times } y_i \neq D(i))/n \tag{F.2}$$

where n is the number of offenders. The AER will usually underestimate the true error rate of the discrimination rule based on the logistic regression. Efron's estimate of the bias in (F.2) is given by:

$$w(p) = 2n^{-1}\Sigma_i p(x_i)[1 - p(x_i)] \, \varphi \, (c_i d_i^{-\frac{1}{2}})d_i^{\frac{1}{2}} \tag{F.3}$$

where

$$c_i = \log[c/(1 - c)] - x_i'b ,$$

$$\varphi(t) = (2\Pi)^{-\frac{1}{2}} \exp (-t^2/2) ,$$

for $t = c_i d_i^{-\frac{1}{2}}$ in this case,

b is vector of fitted logistic regression coefficients, and $d_i$ is the estimated variance of the discriminant score $x_i'b$ and is a quantity available in the output of most logistic regression programs. Specifically

$$d_i = x_i'T^{-1}x_i \tag{F.4}$$

where

$$T = \Sigma_i p(x_i)[1 - p(x_i)]x_i x_i'$$

is the inverse of the usual estimate of the covariance matrix of b. For completeness note that

$$p(x_i) = [1 + \exp(-x_i'b)]^{-1} \tag{F.5}$$

so that

$$x_i'b = \log[p(x_i)/(1 - p(x_i))] .$$

# BIBLIOGRAPHY

Ball, John C., John W. Shaffer, and David N. Nurco, "The Day-to-Day Criminality of Heroin Addicts in Baltimore—A Study in the Continuity of Offence Rates," *Drug and Alcohol Dependence*, Vol. 12, 1983, pp. 119–142.

Belsley, David A., Edwin Kuh, and Roy E. Welsch, *Regression Diagnostics, Identifying Influential Data and Sources of Collinearity*, John Wiley & Sons, New York, 1980.

Bishop, Yvonne, Stephen E. Fienberg, and Paul W. Holland, *Discrete Multivariate Analysis: Theory and Practice*, The MIT Press, Cambridge, 1975.

Blackmore, John, and Jane Welsh, "Selective Incapacitation: Sentencing According to Risk," *Crime and Delinquency*, Vol. 29, 1983, pp. 504–528.

Blumstein, A., and J. Cohen, *Characterization of Criminal Career Patterns from Longitudinal Analysis of Arrest Histories*, School of Urban and Public Affairs, Carnegie-Mellon University, Pittsburgh, 1980.

Blumstein, Alfred, Jacqueline Cohen, Jeffrey Roth, and Christy Visher (eds.), *Criminal Careers and "Career Criminals,"* National Research Council, National Academy of Sciences, Washington, D.C., 1986.

Chaiken, Jan M., and Marcia R. Chaiken, *Varieties of Criminal Behavior*, The RAND Corporation, R-2814-NIJ, August 1982.

Chaiken, Jan M., and Marcia R. Chaiken with Joyce Peterson, *Varieties of Criminal Behavior: Summary and Policy Implications*, The RAND Corporation, R-2814/1-NIJ, August 1982.

Chaiken, Marcia, and Jan Chaiken, "Offender Types and Public Policy," *Crime and Delinquency*, Vol. 30, No. 2, 1984, pp. 195–226.

Chaiken, Marcia R., and Jan M. Chaiken, *Who Gets Caught Doing Crime?* Bureau of Justice Statistics Discussion Paper, Washington, D.C., 1985.

Chelimsky, E., and J. Dahmann, *Career Criminal Program National Evaluation: Final Report*, National Institute of Justice, Washington, D.C., 1981.

Cohen, Jacqueline, "Incapacitation as a Strategy for Crime Control: Possibilities and Pitfalls," in Michael Tonry and Norval Morris (eds.), *Crime and Justice: An Annual Review of Research*, Vol. 5, University of Chicago Press, 1983.

Cook, R. Dennis, and Sanford Weisberg, *Residuals and Influence in Regression*, Chapman and Hall, New York, 1982.

Copas, J. B., "Plotting p Against x," *Applied Statistics*, Vol. 32, No. 1, 1983, pp. 25–32.

Copas, John B., and Roger Tarling, "Some Methodological Issues in Making Predictions," in Blumstein et al., 1986, Vol. 2.

Ebener, Patricia A., *Codebook for Self-Report Data from the 1978 RAND Survey of Prison and Jail Inmates*, The RAND Corporation, N-2016-NIJ, June 1983.

Efron, Bradley, "How Biased Is the Apparent Error Rate of a Logistic Regression?" *Journal of the American Statistical Association*, Vol. 81, 1986, pp. 461–470.

Elliott, Delbert S., Suzanne E. Ageton, and David Huizinga, *The National Youth Survey: 1976 Self-Reported Delinquency Estimates by Sex, Race, Class, and Age*, Behavioral Research Institute, Boulder, 1980.

Fienberg, Stephen E., *The Analysis of Cross-Classified Categorical Data*, The MIT Press, Cambridge, 1977.

Fischer, Daryl R., *Prediction and Incapacitation: Issues and Answers*, presented before the session: "Selective Incapacitation: Methodology and Policy," Annual Meeting of the American Society of Criminology, Cincinnati, November 10, 1984a.

Fischer, Daryl R., *Risk Assessment: Sentencing Based on Probabilities*, Statistical Analysis Center, Iowa Office for Planning and Programming, Des Moines, April 1984b.

Flinn, Christopher, "Dynamic Models of Criminal Careers," in Blumstein et al., 1986, Vol. 2.

Forst, Brian, *The Prosecutor's Case Selection Problem: "Career Criminals" and Other Concerns*, INSLAW, Washington, D.C., 1982.

Forst, Brian, "Selective Incapacitation: An Idea Whose Time Has Come?" *Federal Probation*, Vol. 46, 1983, pp. 10–20.

Goldstein, Paul J., *Habitual Criminal Activity and Patterns of Drug Use*, presented at the Annual Meetings of the American Society of Criminology, November 1982.

Gottfredson, Stephen D., and Don M. Gottfredson, "Accuracy of Prediction Models," in Blumstein et al., 1986, Vol. 2.

Greenwood, Peter W., with Allan Abrahamse, *Selective Incapacitation*, The RAND Corporation, R-2815-NIJ, 1982.

Greenwood, Peter W., and Susan Turner, *Selective Incapacitation Revisited: Why the High-Rate Offenders Are Hard to Predict*, The RAND Corporation, R-3397-NIJ, March 1987.

Hoffman, Peter B., "Screening for Risk: A Revised Salient Factor Score," *Journal of Criminal Justice*, Vol. 11, 1983, pp. 539–547.

Honig, Paul, *Description of Official Record Data from the 1978 RAND Inmate Survey*, The RAND Corporation, N-2017-NIJ, June 1983.

Johnson, Bruce D., *The Drug-Crime Nexus (1979–1981): Research on the Drug-Crime Relationship with Emphasis upon Heroin Users/Injectors as Criminal Recidivists*, prepared for the National Institute of Justice, Washington, D.C., December 1981.

Klein, Stephen P., and Michael N. Caggiano, *The Prevalence, Predictability and Policy Implications of Recidivism*, The RAND Corporation, R-3413-BJS, August 1986.

Landwehr, James M., Daryl Pregibon, and Anne C. Shoemaker, "Graphical Methods for Assessing Logistic Regression Models," *JASA*, Vol. 79, No. 385, March 1984, pp. 61–83.

Lehoczky, John P., "Random Parameter Stochastic-Process Models of Criminal Careers," in Blumstein et al., 1986, Vol. 2.

Marquis, K. H., with P. Ebener, *Quality of Prisoner Self-Reports: Arrest and Conviction Response Errors*, The RAND Corporation, R-2637-DOJ, March 1981.

Petersilia, Joan, "Criminal Career Research: A Review of Recent Evidence," in N. Morris and M. Tonry (eds.), *Crime and Justice*, University of Chicago Press, 1980.

Petersilia, Joan, and Peter W. Greenwood, with Marvin Lavin, *Criminal Careers of Habitual Felons*, The RAND Corporation, R-2144-DOJ, August 1977.

Petersilia, J. R., and M. M. Lavin, *Targeting Career Criminals: A Developing Criminal Justice Strategy*, The RAND Corporation, P-6173, August 1978.

Peterson, Mark A., and Harriet B. Braiker, with Suzanne M. Polich, *Who Commits Crimes: A Survey of Prison Inmates*, Oelgeschlager, Gunn, and Hain, Publishers, Inc., Cambridge, 1981.

Peterson, Mark, Jan Chaiken, Pat Ebener, and Paul Honig, *Survey of Jail and Prison Inmates: Background and Method*, The RAND Corporation, N-1635-NIJ, 1982.

Peterson, Mark A., Suzanne Polich, and Jan Michael Chaiken, *Data Tape for the RAND 1976 Survey of California Prison Inmates*, The RAND Corporation, N-1505-DOJ, June 1980.

Pregibon, Daryl, "Logistic Regression Diagnostics," *The Annals of Statistics*, Vol. 9, No. 4, 1981, pp. 705–724.

Reboussin, David, *Some Diagnostics for Logistic Regression Using the Propensity Score*, Master's Thesis, Department of Statistics, University of Chicago, June 8, 1984.

Reiss, Albert J., *Surveys of Self-Reported Delicts*, unpublished paper, Department of Sociology, Yale University, 1973.

Rhodes, William, Herbert Tyson, James Weekly, Catherine Conly, and Gustave Powell, *Developing Criteria for Identifying Career Criminals*, The Institute for Law and Social Research, Washington, D.C., 1982.

Rolph, John E., Jan M. Chaiken, and Robert L. Houchens, *Methods for Estimating Crime Rates of Individuals*, The RAND Corporation, R-2730-NIJ, March 1981.

Rosenbaum, P. R., and D. B. Rubin, "Assessing Sensitivity to an Unobserved Binary Covariate in an Observational Study with Binary Outcome," *Journal of the Royal Statistical Society*, Series B, Vol. 45, 1983a, pp. 212–218.

Rosenbaum, P. R., and D. B. Rubin, "The Central Role of the Propensity Score in Observational Studies," *Biometrika*, Vol. 70, 1983b, pp. 41–55.

Shannon, Lyle W., "A Longitudinal Study of Delinquency and Crime," in Charles Welford (ed.), *Quantitative Studies in Criminology*, Sage Publications, Beverly Hills, 1978.

Spelman, William, *The Depth of a Dangerous Temptation: Another Look at Selective Incapacitation*, prepared for The National Institute of Justice, U.S. Department of Justice, Police Executive Research Forum, Washington, D.C., February 1986.

Tracy, Paul, Marvin Wolfgang, and Robert Figlio, *Delinquency in Two Birth Cohorts*, Center for Studies in Criminology and Criminal Law, University of Pennsylvania, May 1985.

Visher, Christy, "The RAND Second Inmate Survey: A Reanalysis," in Blumstein et al., 1986, Vol. 2.

von Hirsch, Andrew, "The Ethics of Selective Incapacitation: Observations on the Contemporary Debate," *Crime and Delinquency*, Vol. 30, 1984, pp. 175–194.

von Hirsch, Andrew, *Past or Future Crimes: Deservedness and Dangerousness in the Sentencing of Criminals*, Rutgers University Press, New Brunswick, New Jersey, 1985.

von Hirsch, Andrew, and Donald M. Gottfredson, "Selective Incapacitation: Some Questions About Research Design and Equity," *New York University Review of Law and Social Change*, Vol. 12, 1984, pp. 11–51.

Weisberg, Sanford, *Applied Linear Regression*, 2d ed., John Wiley & Sons, New York, 1985.

Williams, Kristen, *The Scope and Prediction of Recidivism*, The Institute for Law and Social Research, Washington, D.C., 1979.

Wolfgang, Marvin, Robert Figlio, and Thorsten Sellin, *Delinquency in a Birth Cohort*, The University of Chicago Press, 1972.