# ETH

**Swiss Federal Institute of Technology Zurich**

**Seminar for**
**Statistics**

**Department of Mathematics**

Master Thesis

Fall 2018

## Christopher Salahub

# Seen to be done
**A statistical investigation of peremptory challenge**

# Preface

This work would be nowhere near as polished or complete without the effort of Prof. Dr. Marloes Maathuis to ensure I was performing analysis with a clear direction and purpose. I would like to thank her for finding time in her busy schedule to allow for weekly meetings. The group meetings organized by her Ph.D. student Marco Eigenmann were also critical in the development of more nuanced analysis and intuitive visualizations. I thank Marco Eigenmann for organizing them, and Jinzhou Li, Armin Fingerle, Sanzio Monti, and Qikun Xiang for attending my presentations and listening attentively. A special thanks is extended to Cédric Bleutler and Leonard Henckel, both of whom were especially engaged and participated in lengthy discussions both during and outside of the group meetings.

I would like to acknowledge in particular Prof. Dr. Tilman Altwicker for his detailed suggestions on where to look for more legal context on the topic and Prof. Dr. Samuel Baumgartner for his research suggestions. They were very important at providing the necessary information to begin a first investigation of the topic. Of course, without the cooperation of Dr. Ronald Wright, Dr. George Woodworth, Dr. Barbara O'Brien, and Dr. Catherine Grosso for generously providing me with the data from their investigations on the subject of peremptory challenge. Without this data, the visualizations and modelling presented here simply would not have been possible, and so I am exceptionally grateful that they were so enthusiastic to share the fruits of their labour to help cultivate mine.

# Abstract

Short summary of my thesis.

# Contents

# List of Figures

# List of Tables

# Notation and Terms

In order to facilitate clarity despite brevity, a list of terms used in this paper is presented here.

**Prosecution/State**  The legal representation which argues for conviction

**Defence**  The legal reperesentation which argues against conviction

**Court**  Reference to the judge, prosecution, and defence

**Venire**  The population sample from which a jury is selected (according to Mirriam-Webster (2019a) derived from the latin *venire facias*: "may you cause to come")

**Jury**  The final group of (usually) twelve chosen venire members which judge the guilt or innocence of the accused/defendant

**Accused/Defendant**  The individual on trial for a crime

**Voir dire**  From old French "to speak the truth" (see Mirriam-Webster (2019b)), this is the questioning process used by the court to assess the suitability of a venire member to sit on the jury

**Struck**  In the context of a venire member being rejected from the jury, struck indicates removal by peremptory challenge or challenge with cause

**Litigants**  The accusor and the accused

# Chapter 1

# Introduction

The Gerald Stanley murder trial was noteworthy for all of the wrong reasons. The first reason was the crime itself. The rural region around Biggar, Saskatchewan (Quenneville (2018)) is not known for crime, indeed, the crime statistics collected by Statistics Canada suggest it is one of the safest in the province (Statistics Canada (2018)). Any murder at all would be worthy of attention and subject to plenty of drama. But beyond the damage this trial has done to the community, this trial is noteworthy because it led to a significant re-examination of the legal jurisprudence surrounding the jury selection process culminating in the proposition of Bill C-75 by the Canadian government in March of 2018 (42nd Parliament of Canada (2018a)), less than two months after the trial's verdict (Quenneville and Warick (2018)).

Bill C-75, in part, aims to ameliorate one of the critical points of contention about the Gerald Stanley case: the use of peremptory challenges in jury selection. The outsized impact of the case was due, in large part, to it's racial aspect. Gerald Stanley, a white man, was accused of second degree murder in the killing of Colten Boushie, a First Nations man. Given Canada's troubled history with First Nations groups, this alone would have been enough to make the trial a flash point for race issues, but that was not the worst aspect of the trial. Rather, it was the alleged use of peremptory challenges to strike five potential jurors who "appeared" to be First Nations, resulting in an all-white jury, that proved to be the most controversial and influential facet of the entire affair (Harris (2018), MacLean (2018)).

With Bill C-75 currently moving through the Canadian parliamentary system, having completed its second reading in June 2018 (42nd Parliament of Canada (2018b)), a close re-examination of the practice of peremptory challenge is warranted. A great deal of ink has already been spilled on both sides of the debate (see Hasan (2018), Zinchuk (2018), and Roach (2018)), but startlingly little of this discussion has been based on any hard evidence on the impact of peremptory challenge in jury selection. This paper aims to provide analysis and evidence to illuminate the topic further by analyzing three separate peremptory challenge data sets collected in the United States, namely Wright, Chavis, and Parks (2018), Grosso and O'Brien (2012), and Baldus, Woodworth, Zuckerman, and Weiner (2001). While this data cannot tell us if challenges were racially motivated in the Stanley trial, stepping back from this fraught legal episode to take a wider view of the practice of peremptory challenge provides a more sober place to start the discussion of its place in modern jury trials.

This paper will proceed in five parts. Chapter 2 provides a brief history of the practice of peremptory challenges in jury trials, in particular explaining their original motivation, past implementations, and how they have developed in the United States, the United Kingdom, and Canada. Chapter 3 proceeds to discuss the three data sets obtained, explaining the sources and collection methods before detailing the cleaning and preprocessing. Chapter 4 then provides the details and results of the analysis performed on the different data sets. It begins discussing the Jury Sunshine data set, which was used as a 'test' set of sorts, where analysis could be flexibly performed before the final analysis methods were turned to the other two data sets. The results of this analysis are compared to previous works in Chapter **??**. Finally, the results and findings are summarized in **??**, and recommendations based on the observations obtained here are provided.

# Chapter 2

# Peremptory Challenges

Although the focus of this text is the legal practice of peremptory challenges, these are a specific practice which may not be known in detail to the reader, a brief exploration of their history, motivation, and current use is presented here. It is not meant to be exhaustive, but rather to provide context and references for an interested and motivated reader to learn more. Indeed, many details have been omitted from the summary of the history in particular.

## 2.1  Jury Selection Procedures

Before reviewing the history, it is best to give some context and an explanation for readers unfamiliar with the jury system and general courtroom procedures. The general steps shared by jury trials are outlined below. More detail and a discussion of the diversity of jury selection procedures can be found in Ford (2010), Hans and Vidmar (1986), and Van Dyke (1977).

i.) Eligible individuals are selected at random from the population (using a list known as the *jury roll*) of the region surrounding the location of the crime, the sampled individuals are called the *venire*

ii.) The venire is presented to the court, either as a group or sequentially (borrowing the names of Ford (2010): the "struck-jury" system and the "sequential-selection" system, respectively)

iii.) The presented venire member(s) are questioned in a process called *voir dire*, which can result in three possible outcomes for each venire member:

   (a) The venire member is removed with cause, the cause provided by either the prosecutor or defence lawyer and admitted by the judge

   (b) The venire member is removed by a peremptory challenge by the prosecutor or defence lawyer, where no reason need be provided to the court

   (c) The venire member is accepted into the jury, and so becomes a juror

iv.) Steps i-iii are repeated until the desired number of jurors has been found

The details of each of these steps varies by region. Jury rolls can be collected from many different sources. In the United States, they are typically selected using lists of registered voters (see Van Dyke (1977) chapter two), but in Canada the practice varies province by province. Ontario uses a combination of municipal voter lists and First Nations band lists (see Ministry of the Attorney General of Ontario (2018)), while in Saskatchewan - the province of the Gerald Stanley trial - the jury roll is created from the data in the central government health insurance agency in accordance with Government of Saskatchewan (1998).

While two presentation methods are observed in step ii, Ford (2010) and Van Dyke (1977) both note that the predominant method in the United States and Canada is the sequential-selection system. This is perhaps due to the relative efficiency of the method, as it is clear that in the sequential system voir dire need not be performed on the entire venire, only a subset. Contrast this with the struck-jury system, where the entire venire must be reviewed in every trial.

Finally, the scope of voir dire is radically different in the United States and much of the British Commonwealth. Van Dyke (1977) notes that Canada and the United Kingdom do not allow questions in areas of "non-specific" bias, or bias which is not directly related to the case before the court. That is to say, while it would be perfectly valid to ask a venire member for a murder case about their work history in the United States, such a question would only be allowed in Canada or the United Kingdom if occupation was specifically related to the murder case. Hans and Vidmar (1986) suggest that this difference is due to a difference in philosophy. To borrow a quote from page 63 of Hans and Vidmar (1986):

> In Canada, for example, the courts have said that we must start with an
> initial presumption that "a juror will perform his duties in accordance with
> his oath"

This opinion places a greater responsibility on the jurors themselves to overcome their biases and accept arguments in spite of them. The American opinion that certain prejudice cannot be overcome by jurors stands in stark contrast.

## 2.2   The Role of the Jury

The central function of a jury in a jury trial system is to judge the innocence or guilt of an accused in light of evidence. This has varied drastically in form throughout history. Consider that in the distant past, von Moschzisker (1921) and Hoffman (1997) report that the central function of the jury was to collect evidence, essentially assuming the role commonly performed today by police detectives. Such a role justified the practice of selecting the most "trustworthy" individuals of some reknown.

This is contrasted by the modern jury, which performs no collection of evidence, and is meant to be composed of a panel of peers or "equals" of the accused sampled at random from the population, an idea which did not develop until 19th century Britain (see page 28 of Hans and Vidmar (1986)) and was not applied using random sampling until some time later (see Hoffman (1997) and page 29 of Hans and Vidmar (1986)). The modern jury is meant to apply the law, as told to them by the judge[1], to the case at hand. Evidence for guilt is then presented to the jury by the prosecutor, while evidence meant to exonerate is presented by the defence.

The jury listens to the evidence, considers the law as presented by the judge, and must (typically) reach a unanimous decision of guilt or acquittal. Such a decision cannot be overturned by the judge of the court, and the judge must then determine sentencing based on the decision of the jury and the letter of the law[1]. It should be clear that the jury therefore has tremendous power in the legal system. The philosophical and ethical justification for such power is well explained by Woolley (2018), and best summarized by a quote from Supreme Court of Canada (1991):

> The jury, through its collective decision making, is an excellent fact finder; due to its representative character, it acts as the conscience of the community; the jury can act as the final bulwark against oppressive laws or their enforcement; it provides a means whereby the public increases its knowledge of the criminal justice system and it increases, through the involvement of the public, societal trust in the system as a whole.

While such enthusiastic support for juries has not been expressed by all countries which practice them, the justification is entirely consistent with the histories and discussions presented by Hoffman (1997), von Moschzisker (1921), Hans and Vidmar (1986), Van Dyke (1977), and others.

## 2.3   Modern Peremptory Challenge Controversy

If the general utility and importance of the jury is clear, the same cannot be said for peremptory challenges. The privileged privileged removal of a venire member - to be replaced by a new randomly selected venire member - by either the prosecution or defence without justification has seen allegations of abuse.

In the United States, repeated allegations of racial discrimination have led to significant changes in their allowed use, through cases such as Supreme Court of the United States (1965) and Supreme Court of the United States (1986). The first of these cases, *Swain v. Alabama*, established in 1965 that the systematic removal of venire members of a particular race could be unconstitutional discrimination under the Fourteenth Amendment, but argued that a *"prima facie"* (or "based on first impression") argument of discrimination was not adequate to prove this. This placed a significant burden on the side taking issue with a challenge to demonstrate discrimination in the use of peremptory challenges.

However, this ruling was overturned only 21 years later in the 1986 case *Batson v. Kentucky*, which allowed the party ojecting to a challenge to use a *prima facie* argument which must be countered by a race-neutral reason that satisfies the judge. If no such reason can be supplied, the challenge would not be allowed. This created a new challenge to use against peremptory challenge to keep a venire member: the so-called "Batson Challenge". While the effectiveness of this system of additional challenges is questionable both practically and in abstract (see Page (2005) and Morehead (1994), and a particularly strong response in Hoffman (1997)), it has only been extended to allow Batson challenges for gender and other characteristics of venire members.

In Canada, there have also been racial controversies. A report by a government inquiry in the province of Manitoba in 1991 (see Roach (2018)) was already reporting on possible

---

[1]Hans and Vidmar (1986) note that this system actually varies throughout the US, though the jury and judge powers described here are consistent across Canada.

racial bias against First Nations venire members. More damning still was the Iacobucci (2013) Report on First Nations representation in juries proposed an explicit restriction to the practice when it recommended:

> ...an amendment to the Criminal Code that would prevent the use of peremptory challenges to discriminate against First Nations people serving on juries.

Despite these recommendations and allegations, there had not been a significant political effort to reform the peremptory challenge system until the Gerald Stanley trial culminated in the tabling of Bill C75 42nd Parliament of Canada (2018b). As it currently stands, the bill has not been approved by the Government of Canada, but seems likely to become law in the near future, which would abolish the peremptory challenge in Canada.

In doing so Canada would join the United Kingdom. Significant controversy around the use peremptory challenges there already led to the abolition of the practice by parliament in the Criminal Justice Act of 1988. The specific controversy was the result of the Cyprus spy case in the late 1970s, which led to a "sustained campaign in Parliament and in the press alleging that defence counsel were systematically abusing it" (see Hoffman (1997))[2].

## 2.4 The Role of the Peremptory Challenge

Despite these legal changes, recommendations, and a great deal of articles providing analysis against the practice (see, for example, Hoffman (1997)), the topic remains controversial. The modern motivation and justification for the practice in spite of all of the controversy is perhaps best described by Justice Byron R. White in Supreme Court of the United States (1965):

> The function of the challenge is not only to eliminate extremes of partiality on both sides, but to assure the parties that the jurors before whom they try the case will decide on the basis of the evidence placed before them, and not otherwise. In this way, the peremptory satisfies the rule that, "to perform its high function in the best way, justice must satisfy the appearance of justice."

Such a justification harks back to the now famous words of Lord Chief Justice Hewart in *R v. Sussex Justices* in 1924: "Justice should not only be done, but should manifestly and undoubtedly be seen to be done" (as reported in Richardson Oakes and Davies (2016)). While these words originally only referred to the pecuniary interest of court staff involved in the case, they have since come to express the idealized expectation that both the defence and prosecution find the judge and jury acceptable, as explored by Richardson Oakes and Davies (2016)[3].

This defence suggests two modern justifications for the peremptory challenge. The first is that of removing venire members with "extreme" bias, and the second is the creation

---

[2]It should be noted that this did not abolish the use of "standing-aside" by the Crown, although the practice has been heavily curtailed to only national security trials with strict guidelines to its use, which are outlined by Attorney General's Office of the United Kingdom (2012).

[3]Such grand generalizations and myth-making can also be seen in the common belief that the right to a trial by jury was originally established in the Magna Carta, an idea which is not supported by the relevant historical evidence (see Hoffman (1997) and Van Dyke (1977) for a detailed discussion and more accurate history).

of a jury which is composed of jurors mutually acceptable to both the defense and the prosecution. Those who defended the practice of peremptory challenges in Canada after the Gerald Stanley trial, including Hasan (2018) and Macnab (2018), seem to use this defence or some variant of it to argue in favour of keeping the practice. However philosophically appealing these two claims are, in light of all of the controversy surrounding the peremptory challenge, perhaps a critical and empirical examination of these assertions is warranted.

## 2.5 History

Such an analysis might appropriately begin with a historical explanation of the peremptory challenge. Roughly, the presentation of the history of jury trials here follows the comprehensive and exhaustively referenced description provided by Hoffman (1997). Two of the references Hoffman uses extensively, Hans and Vidmar (1986) and Van Dyke (1977), provided useful context while specific details provided by von Moschzisker (1921), Forsyth (1994), Brown, McGuire, and Winters (1978), and Brown (2000) helped to create a clearer picture of particular periods of jury history. Information regarding the history of the Canadian system was provided by Brown (2000) and Petersen (1993). For an excellent exploration of the nineteenth century, a formative time for the development of challenge law, see Brown (2000).

It must be noted that certain important trials in the development of the peremptory challenge system have been excluded from the summary provided here. This was done deliberately, as the history presented here is only meant to present the practice of peremptory challenges throughout history in broad terms. All of the sources listed above are much more thorough, by merit of their singular focus on the analysis of the practice from a legal and historical perspective, while this work devotes more to empirical and statistical analysis.

### 2.5.1 Pre-English History

Although precise timelines are hard to establish, there is evidence that jury trials have occurred in some form or another since antiquity. The concept, that of judgement by a group of peers, is so ancient that it is prevalent not only in historical records, but in myth. As Hoffman (1997) indicates, both Norse and Greek mythology feature groups of individuals assessing the guilt or collecting evidence about the actions of a peer.

Outside of the realm of myth, Hoffman (1997) reports on evidence of the use of juries in Ancient Egypt, Mycenae, Druid England, Greece, Rome, Viking Scandanavia, the Holy Roman Empire, and Saracen Jerusalem. It should be noted that in none of these areas was the jury trial the primary form of conflict resolution practiced. Nonetheless, it is clear the jury trial has a broad and long history of use.

Something similar to the modern peremptory challenge does not appear until Rome, however. The Roman *Judices* were groups of senators selected to judge the guilt of the accused in a legal case. Hoffman (1997) presents evidence of the selection of 81 Senators to sit on one of these *Judices*, after which the litigants were permitted to remove fifteen of these

Senators each. This egalitarian reduction of the jury size seems analogous to the modern peremptory challenge system in placing the power of removal with the litigant and suggesting no justification is necessary for their removal.

### 2.5.2  In English Law (1066–1988)

Peremptory challenge did not reach is modern form, as outlined above, until it was established in the English legal system. It should be noted that despite some previous debate on the topic, the most modern historical evidence suggests that the basis of the English practice was not related to the system used in the selection of *Judices* in Rome. The English system appears to be its own beast entirely.

The dominant historical interpretation is presented by von Moschzisker (1921) and Hoffman (1997): that the jury system was introduced to England during the Norman conquest of 1066 by William the Conqueror. The practice, however, was not made official until the Assize of Clarendon in 1166 by Henry II, and it was not until the outlaw of trials by ordeal (the most common method of trial at that time) in 1215, that peremptory challenges began to appear in England in the late thirteenth century. The challenges were officially recognized in 1305 when Parliament outlawed their use by the Crown, only to replace them with an analogous system of so-called "standing-aside"[4].

It should be noted here that although the challenges issued between the Assize of Clarendon and this 1305 act are called "peremptory," they may not have served the same purpose, nor the same justification, as the modern challenges. Indeed, as Hoffman (1997) argues convincingly, these challenges may have been closer to modern challenges with cause. The argument hinges on the paradigm of royal infallibility and absolutism which was present in the late medieval period when the peremptory challenge first appeared (see Burgess (1992)).

Under royal absolutism and infallibility the argument for peremptory challenges is quite simple. If the king cannot be wrong in his judgement and he has some reason to feel that a venire member cannot serve on the jury, then he need not say why he thinks that is so, as his judgement is correct in any case. Indeed, asking for an explanation would be disrespectful and providing one undignified. The Crown prosecutors, as representatives of the king, would be similarly shielded from criticism.

Additionally, this is supported by the abolition of their royal use in 1305, the language of which suggests that peremptory challenges were originally the privilege of the Crown (see Hoffman (1997) and Van Dyke (1977)), with none being granted to the defence. As royal infallibilty grew out of favour, peremptory challenges seem to have been granted to the defence, rather than being removed entirely.

Whatever the logic of the expansion of these challenges to the defence, their legal limits are recorded more precisely[5]. From a maximum of 35 challenges allowed at their peak in the fourteenth century, the number of challenges allowed only decreased over time until their abolition in 1988 (discussed in 2.4).

---

[4]For a detailed explanation of this system see Hoffman (1997) and Brown (2000)

[5]see Brown (2000) for a detailed examination of the case law developing around challenges in the nineteenth century

### 2.5.3   In American Law (ca. 1700–1986)

von Moschzisker (1921), Hoffman (1997), and Van Dyke (1977) all agree that the early English colonists that came to North America accepted the jury system with peremptory challenges as common law well before the establishment of the United States of America. Hans and Vidmar (1986) note, however, that the difficulty of ocean travel and the overall indifference of appointed Crown representatives in the colonies led to an increased importance of the jury trial and the role of challenges to these early colonists as a way to exercise some degree of community control in the face of laws drafted in a distant country and implemented by unsympathetic authorities[6].

It is somewhat interesting then, that the United States constitution makes no mention of the practice of peremptory challenges. The Sixth and Seventh Amendments specify a great deal of the jury system, including the right to public defense and an impartial jury drawn from the district of the crime, but make no mention of a right to the exercise of peremptory challenges, or any challenges whatsoever (see Constitution of the United States (1788)).

As Hans and Vidmar (1986) report on page 37, an original draft of the Sixth Amendment expressly included challenges for cause, but the debate around their inclusion resulted in the removal of their mention. They continue to say that at the time, even some proponents of the challenge considered the reference unnecessary, as the practice was implied by the text which remained, referring to a trial by an "impartial" jury. Another result of these debates was the adoption of the extensive voir dire process which allows questions of general bias (page 37-38 of Hans and Vidmar (1986), though Brown (2000) notes that 1807 Burr trial was also highly significant in the development of general voire dire in the United States).

Critically, there appears to have been no discussion around the inclusion of peremptory challenges (see Hans and Vidmar (1986) and Hoffman (1997)). Despite the clear importance of the jury trial to the drafters of these amendments, it would seem the peremptory challenge was not considered to have anywhere near the same significance as judgement by an impartial jury of local peers[7].

Regardless of this, as Brown (2000) notes, the importance and use of challenges increased in the United States in the nineteenth century following American independence due to a desire to prevent the tyranny of the state. This desire also led to the adoption of a limited number peremptory challenges for the prosecution, rather than the possibly unlimited stand-asides that were allowed under British law to prosecutors (see Van Dyke (1977), page 150).

While the specific numbers of peremptory challenges allowed to both sides and the required motivation of challenges for cause have varied over time (see Hoffman (1997) and Brown (2000)), they have remained a feature of the American legal system, and numerous

---

[6]For more detail on this development among the early colonists, it is instructive to read about the Zenger trial of 1734 (described form pages 33-35 of Hans and Vidmar (1986)). Not only does this trial say a great deal about the attitudes of the colonists at the time, but it also presents the idea of a jury assessing guilt and "wrongness" using their own conscience rather than just settling fact, the precept of the modern jury trial in Canada (see Woolley (2018)) is based on this very idea

[7]Indeed, as *Batson v. Kentucky* and *Swain v. Alabama* have both shown (Supreme Court of the United States (1986) and Supreme Court of the United States (1965)), the modern interpretation of "impartial" may preclude the use of peremptory challenges altogether

Supreme court cases (detailed by Hoffman (1997)) have merely served to make the use of challenges more specific and codified. It was not until *Batson v. Kentucky* in 1986 that this system of challenges was drastically changed with the introduction of Batson challenges which would nullify peremptory ones.

### 2.5.4   In Canadian Law (ca 1800–2018)

Canadian law, inspired by a close relationship to both the British Crown and the United States, seems to have adopted elements of both legal systems in its development of peremptory challenges in the nineteenth century. As discussed by Brown (2000), Canada adopted the American practice of replacing prosecutorial stand-asides in favour of the more egalitarian granting of limited peremptory challenges to both sides. Despite this, the Canadian voir dire process remains limited and much more similar to the British one, as does the system of challenges for cause (see page 48 of Hans and Vidmar (1986)).

One perfect demonstration of this departure is the Canadian constitution. As in the United States, the Canadian consitution fails to mention challenges. The British North America Act of 1867 (see Constitution of Canada (1982)), which established Canada's independence from England, makes no mention of legal rights of the accused, indicating a deference to legal precedent as in England. It is not until the Charter of Rights and Freedoms in 1982[8] that such rights were guaranteed in a founding document. Notably, its language is considerably more vague than the United States Sixth and Seventh Amendments, guaranteeing only "the benefit of trial by jury" (see Constitution of Canada (1982)).

This "eclectic" incorporation of both American and English case law, to borrow the term used by Brown (2000), led to a system somewhere between the English and American systems, but decidedly closer in operation to the English system. It should be noted, however, that as Canada grew more populous in the twentieth century and developed a greater legal precedent and more experienced judges of its own, this reliance upon its former colonial master and its more powerful southern neighbour seems to have diminished in importance. Viewing Supreme court rulings from recent decades reveals a great deal of reference to Canadian legal precedent rather than to the precedent of the other two countries.

## 2.6   Summary

The peremptory challenge has existed in some form since at least the fourteenth century. After its inception in England, it spread with English conquest and colonization, with new colonies and local governments accepting the practice based primarily on the adoption of English legal precedent. Despite its abandonment by the English in 1988, it has remained highly prominent in the United States, accompanied by a voir dire process far more thorough and invasive than that performed in the English or Canadian jury selection process.

---

[8]This was the year of patriation of the Canadian constitution. As independence was granted by the Britsh Parliament, the British North America Act outlining Canada's laws was a British law and changing it was the prerogative of the British Parliament rather than the Canadian one. It was not until the Consitution Act of 1982 that the Canadian constitution became a Canadian law. For a more detailed history see Sheppard (2018)

Though the practice has historical longevity, it is not guaranteed by the constitutions of Canada or the United States, and has been a practice of considerable legal debate and significant change throughout its history. In England this culminated in the Cyprus sky trial, in the United States in *Batson v. Kentucky* and *Swain v. Alabama*, and in Canada in *R. v. Stanley*: the Gerald Stanley murder trial. As a consequence, the broad agreement of the importance and propriety of a jury has conferred little consensus on the place peremptory challenges in the selection of juries.

# Chapter 3

# Data

Without data, performing an analysis that incorporated more than the history and legal argumentation presented in Chapter 2 is impossible. This proved problematic. While the motivation of this text was a Canadian case, no comprehensive Canadian data sets which exmained jury selection in Canada could be found. The increased prominence of the jury selection process in the United States garnered a more fruitful search.

The author is heavily indebted to Wright et al.; Grosso and O'Brien; and Baldus et al.. These authors shared their data freely with the author, providing him with a wealth of data to analyse empirically. As a consequence of the multiple separate data sets, however, care must be taken to describe each of the data sets separately in order to capture adequately the different methodologies and sources they represent. Critically, it should be noted that each of these papers represents effort on the part of the authors. As Wright et al. (2018) notes:

> limited public access to court data reinforces the single-case focus of the legal
> doctrines related to jury selection. Poor access to records is the single largest
> reason why jury selection cannot break out of the litigato's framework to
> become a normal topic for political debate

Currently, the collection of jury data is difficult, as many courtrooms have not digitized past records and concerns over privacy limit the release of those records, which are stored as paper documents in the case file (see **?**). This limits the ability of an individual to ask for summaries across numerous trials or to view the jury selection process on a scale beyond the basis of one case. Thus, to gather aggregate data the authors of these papers necessarily used different collection techniques dictated by the scope of collection desired and the procedures of the court systems from which data was collected.

## 3.1 Jury Sunshine Project

### 3.1.1 Methodology

The Jury Sunshine Project (Wright et al. (2018)), so named as it was carried out in order to shed light on the jury selection process, is the most extensive data set which was provided to the author. It endeavoured to collect jury data for all felony trial cases in

North Carolina in the year 2011, which ultimately resulted in a data set that detailed the simple demographic characteristics and trial information of 29,624 individuals summoned for jury duty in 1,306 trials. Note that not all entries were complete.

Due to the scope of the project, there are a number of problems which had to be solved by the authors. The first of these was simply identifying which court cases went to trial in 2011, in order to direct resources effectively. This was accomplished by downloading publicly available case data from the North Carolina Administrative Office of the Courts (NCAOC)[1] and determining the case numbers and counties of cases which went to trial. Wright et al. state that this likely missed some cases which went to trial, but that they were confident that a "strong majority" of trials was collected, which did not systematically differ from those excluded.

This list was then used to perform a pilot study to refine recording practices before undertaking a more general survey where "law students, law librarians, and undergraduate students" (called *collectors* for convenience) visited court clerk offices to collect the relevant case data, including the presiding judge, prosecutor, defence lawyer, defendant, venire members, charges, verdict, and sentence. The case files also included data about whether a venire member was removed by cause or peremptorily, and the party which challenged in the peremptory case. Using public voter databases, bar admission records, and judge appointment records, these collectors were able to determine demographic (race, gender, and date of birth) and political affiliation data for the venire members, lawyers, defendants, and judges. This data set was stored stored in a relational database provided to the author by Dr. Ronald Wright.

The analysis of the data provided in Wright et al. (2018) was limited to aggregate summaries of the trends at the venire member level. That is to say, they examined the strike trends for both the defence and the prosecution, conditioning on some additional variables. There was also spatial analysis performed, were different urban counties were directly compared. These analyses were not statistical in nature, and were displayed using contingency tables. Regardless the stark differences between prosecution and defence with regards to race were a key finding when the aggregate data was analyzed.

### 3.1.2 Cleaning

**Flattening the Data**

For greater expediency of analysis, the relational database of the Jury Sunshine Data was first flattened. The relational database was read into Microsoft Excel and the `readxl` package (Wickham and Bryan (2018)) was used to read the excel file into the programming language R . A wrapper for the `merge` function was developed which provided simple a simple output detailing failed matches in an outer join in order to ensure that the flattening of the data into a matrix did not miss important data due to partial incompleteness. The code for this wrapper can be seen in A.2.

---

[1]The link provided in the Jury Sunshine Paper to the specific source (http://www.nccourts.org/Citizens/SRPlanning/Statistics/CAReports_fy16-17.asp) does not appear to be working as of January 2019, however the NCAOC seems to provide an API functionality at https://data.nccourts.gov/api/v1/console/datasets/1.0/search/

This wrapper revealed only a small number of irregularities in the data, which are detailed in A.3:

i.) Twenty-nine charges missing trial information such as the presiding judge (all of trials with IDs of the form 710-0XX)

ii.) Twenty-six prosecutors not associated with any trials and missing demographic data

iii.) One trial missing charge information

Ultimately, the jurors for trial ID 710-01, the trial missing a charge from above, were included in the data as their records were complete otherwise. The prosecutors and charges which could not be joined were excluded from any future analysis, as they could have easily been included by collectors by accident. Due to the small relative size of these inconsistencies relative to the size of the data set, they did not cause concern.

### Uninformative Columns

Of course there were other irregularities in the data than the obvious ones that arose in the flattening process. There were a handful of likely sources for these errors. The first of these is the anonymization of the data for public use. The private data includes a wealth of privileged data such as juror name and address, and these were removed in the data given to the author.

As a consequence of this anonymization as well as the inclusion of rarely used columns such as those for additional notes, some columns of the data contained only NA values. Most baffling of these was the `BirthDate` variable in the `Jurors` table, as there was no clear reason for this data to be missing. Thankfully, none of the missing columns were relevant to the joins performed in flattening, and they would have been only secondary in data analysis. As a consequence, these uninformative columns were simply removed from the data.

### Coding Inconsistencies

Related to this problem was the issue of inconsistently coded variable levels. An example of these inconsistencies would be levels recorded as both lower and upper case letters, or the presence of "?" instead of "U" for unknown values. It is very likely this inconsistency was a direct result of the data collection method which used many data collectors working independently in different places at different times. Thankfully, Wright et al. provided the codebook used by data collectors, which served as the authoritative reference for the admissible factor levels of demographic variables. Solving this problem was as simple as setting all demographic variable levels to be uppercase and replacing obviously mis-specified levels.

One specific inconsistency which should be noted is that of the outcome, which had a handful of entries recorded as "HC", an inadmissible level not defined by the codebook. It is quite possible that this level represented a typo, as the "H" and "G" keys are adjacent on the American QWERTY keyboard layout, and "GC" was the code for 'guilty as charged', but out of a desire to be conservative with the data the occurrence of this outcome was replaced with U instead. Additionally, the inadmissible level "G" was replaced by GC.

**Swaps**

A more difficult problem of level misspecification was the presence of what appeared to be columns with swapped values, frequently occurring with the gender column (the admissible levels of which are "M", "F", and "U") and the political affiliation column (the admissible levels of which are "D", "R", "I", or "U"). The aforementioned "swaps" appeared as records in which, for example, the gender was recorded as "R" and political affiliation as "M". More complicated swaps of three columns also occurred. To address this problem, the `IdentifySwap` function was written (see line 108 in A.2).

The `IdentifySwap` function accepts two arguments: a data frame with named columns and a named list of vectors of the acceptable levels for some of the column names. It then performs vectorized checks of the specified column names and presents any rows which may have swaps or errors interactively to the user, along with a suggested reorder to "un-swap" the row. The user can press enter to accept the suggested reordering, enter some other reordering, or enter 0 to indicate that the row was not a true swap, but simply an error. The un-swapped entries are then returned to the data, and the rows with errors have the erroneous values replaced by "U", the universal code for "Unknown" within all data variables[2].

The source of these swaps is also most likely the data collection method. The codebook provided specifically notes that the data collection was meant to record the race (R), gender (G), and political affiliation (P) data in the form RGP, but it is not inconceivable that it would occasionally have been recorded or entered in some other ordering in the tedium of data entry. In any case, this problem affected only 431 records of the nearly 30,000, suggesting that the recorded error rate was not unacceptably large.

**Charge Classification**

Perhaps the least regular data in this data set was that of the charge text. Due to the lack of any codebook guidance about the standard way of recording a charge in a trial, identical charges were recorded in numerous ways. The first method used to combat this was removing non-alphnumeric characters, extra spaces, and converting all charges to lower case. This still left a considerable variation, however. Consider the charge of breaking and entering, for example, even with this simple preprocessing performed the entries varied significantly (e.g. "break or enter", "breaking andor entering", "breaking and or entering", etc.).

As a consequence, the processing was more involved. First the most common versions of the charge text for the charges were all regularized to be identical (see `StringReg` in A.2. Next, a regular expression classification tree was developed, which would also account for specific features of a charge. When identifying murder, for example, it seemed important to ensure attempted murder was separated from murder itself, and separating first and second degree was also desired. This tree would, when presented with a charge, apply the regular expressions at each node to the charge. If the charge matched the expression at a node, the regular expressions of that node's children were applied to the charge until it was classified to some leaf node, each of which had a standardized value which replaced

---

[2]One notable exception to this insertion of "Unknown" was the case of the judge Arnold O Jones II, whose gender was not recorded in the data, but who was identifiable as a man using a quick Google search of his unique name

the charge. A small example of this structure is displayed in Figure 3.1, and the full tree is visualized in A in Figure A.1.

```
                                              ┌─────────────┐
                                              │   first|1   │
                                              └─────────────┘
                        ┌─────────────────┐
                        │ sex(?=.offense) │
                        └─────────────────┘
                                              ┌─────────────┐
                                              │  second|2   │
                                              └─────────────┘


                                              ┌─────────────┐
                                              │    regis    │
                                              └─────────────┘
                        ┌──────────────────┐
                        │ sex(?=.*offend)  │
                        └──────────────────┘
                                              ┌─────────────┐
                                              │    addr     │
                                              └─────────────┘
```

Figure 3.1: A example of a simple charge classification tree to separate the sexual offenses from charges leveled against previously known sex offenders. A charge would be classified from most general on the left to most specific on the right.

By performing regularization using this charge tree, regularized charges were guaranteed. The cost of this regularization was the inability to classify all crimes, however. Of the 1407 charges present in the data, the tree provides regularization for 1209. With additional time and inspection of the failed matches, the tree could conceivably be axpanded to regularize all charges. As this was not the primary investigation of the report, however, such effort was not expended.

Instead, a number of helpful aggregation and extraction functions were developed to allow for the charges to be further simplified, as seen in A.2. In this report, they were aggregated by intuitive classes: sex-based offenses, thefts, murders, drug charges, violent offenses not otherwise classified, and driving charges. However, other classes, such as the North Carolina felony classes themselves (as provided by North Carolina Sentencing and Policy Advisory Commission (2017)), may provide a more informative classification rationale.

**Variable Level Renaming**

The final step of the data cleaning process was to convert the uninformative codes used to indicate variable values to more intuitive and clear names (for example to convert "I" in the political affiliation variable to "Ind", a clearer indication of independent). Certain variables which were already clear, such as gender (codes "M", "F", "U"), were not renamed due to the clarity of the one letter representations.

### 3.1.3   Variable Synthesis

In order to expand the possible analysis and visualization potential, a number of variables were synthesized from the Jury Sunshine data set. They are detailed below.

**Race Match**  A logical variable which is true for a venire member if they are the same race as the defendant, and false otherwise. This variable was motivated in particular by the Gerald Stanley trial, in which the implicit contention of those who have taken issue with peremptory challenge is that the First Nations venire members were removed by the defence as their race did not match that of Stanley in the racially charged trial.

**Guilty**  Logical indicator indicating whether the trial outcome was guilty or not

**Visible Minority**  Logical indicator of non-white venire member race

**Race of Striking Party**  Factor variable which gives the race of the prosecution if the venire member was struck by the prosecution, the race of the defence if the venire member was struck by the defence

**Simplified Race**  Due to the scarcity of the other minority races, this variable simplified the race provided to white, black, or other for the venire member

**Simplified Defendant Race**  The same as the simplified race for the defendant races

## 3.2   Stubborn Legacy Data

### 3.2.1   Methodology

Grosso and O'Brien (2012) also provided data to the author, albeit a more limited set. This study, also based in North Carolina, focused on the trials of inmates on death row as of July 1, 2010, yielding a total of 173 cases. In each proceeding, the study examined only those venire members not excluded for cause, and critically the analysis of the study focused only on prosecutorial peremptory challenges. Besides collecting demographic data as in the Jury Sunshine Case, this study also collected attitudinal data for the venire members.

Staff attorneys from the Michigan State University College of Law were responsible for the data collection in this study. The work was performed similarly to the Jury Sunshine Data, using case files to collect information about the court proceedings such as the peremptory challenges used, presiding judge, prosecutor, and defence lawyer. Detailed verdict and charge information was not collected, as the preselection criteria of death row inmates made the verdict clear, and the death penalty can only be applied for certain crimes.

To collect demographic and attitudinal data, the juror questionnaire sheets were consulted[3]. These sheets are typically used as a component of voir dire, in order to make the process more efficient and determine venire members categorically ineligible for jury duty. As a result, they inquire about opinions on the death penalty, for example, as well as

---

[3]As Grosso and O'Brien (2012) observe, self-identified race may be the most accurate source of racial group identification

demographic questions. As not all jury questionnaires were available, additional information was collected from jury roll lists to determine the races of the final jury members. It should be noted that this collection was done blind and to high standards of proof, and a reliability study carried out in Grosso and O'Brien (2012) indicated that under this system the race coding was 97.9% accurate when the standards were met. Those for whom the standards were not met were marked as "Unknown."

The analysis performed in this paper was more statistical than in the Jury Sunshine Data. Contingency tables generated using the data were tested using Chi-squared independence tests, and a simple logistic regression model was created to predict prosecutorial strikes. One minor criticism which could be made of their methodology is the lack of a consistent level to their tests. It seems that rather than class these tests as significant or not, these tests were simply performed to report the p-values they returned. Additionally, there are possible multiple testing issues as the study seems to indicate multiple tests were performed on each table, with the specific test used to generate the reported p-values not clearly indicated.

### 3.2.2 Cleaning

## 3.3 Philadelphia Data

### 3.3.1 Methodology

Baldus et al. (2001) presents a similar data set collected using similar means. Court files such as the juror questionnaire, voter registration, and census data were all used to complete juror demographic information for 317 venires consisting of 14,532 venire members in Philadelphia capital murder cases between 1981 and 1997[4]. It should be noted that this data included only those jurors kept or peremptorily struck, venire members struck for cause were not included in the data. The procedure used to determine race using the census and voter registration polls was quite complicated, but was rigorously performed using accepted census methods to a standard of 98% reliability[5].

In their incredibly thorough analysis of the data, there were findings consistent with both the Jury Sunshine and Stubborn Legacy data. The defence and prosecution seem to follow mirrored patterns of racial preference in the use of peremptory challenges, even when controlling for other possible confouding effects.

### 3.3.2 Cleaning

---

[4]This study took into account the sampling error by reweighting venires based on the year of the trial and the defendant race, as court records showed that the sample coverage varied over these factors

[5]Additionally, imputation was only performed in a small minority of cases

# Chapter 4

# Analysis

With this data cleaned and processed, questions can now be posed and addressed through analysis. A few obvious questions come to mind, considering the previous work done on this subject. The first is whether the results found by previous analyses which did not use statistics are statistically significant. Additionally, we may wonder whether the most common arguments posed in favour of peremptory challenge are satisfied in this data.

## 4.1   Case Level Summary

While Wright et al. (2018) reported a great deal of aggregate statistics about the venire members themselves, one piece of investigation which was lacking was an analysis which aggregated and viewed the trends for the cases, rather than simply for individual venire members. As we cannot know why a potential venire member is struck individually, and viewing their aggregate statistics tells us nothing about how different strikes relate to each other, it is possible we are viewing some effect which is not a result of persistent bias across trials, but is rather the result of some other effect.

By aggregating the venire members by trial and viewing the demographic trends in strikes and behaviour at this level, we gain a more detailed insight into the impact of challenges at a more relevant scale. Additionally, such aggregation allows for the synthesis of certain measures, such as a disitributional difference via the Kullback-Leibler divergence (Kullback and Leibler (1951)), which would otherwise not be well defined. This particular perspective of the data has also not been explored by any other studies known to the author.

## 4.2   Modelling

In order to create a single model to test the statistical significance of the differences observed for strike rates by race, defendant race, and party doing the striking, a saturated poisson regression model was fit to the data. Letting $i$ denote the level of the venire member race, $j$ the defendant race, and $k$ the disposition, the numbers of observed venire members in each $ijk$ combination, $y_{ijk}$ were modelled as Poisson-distributed random variables with expectation $\lambda_{ijk}$. A saturated model was then fit to the data, that is a model described by the equation:

$$\log E[y_{ijk}] = \mathbf{x}_{ijk}\beta = \beta_o + \beta_R x_{i..} + \beta_D x_{.j.} + \beta_S x_{..k} + \beta_{R:D} x_{i..} x_{.j.} + \beta_{R:S} x_{i..} x_{..k} + \beta_{D:S} x_{.j.} x_{..k}$$
$$+ \beta_{R:D:S} x_{i..} x_{.j.} x_{..k} \quad (4.2.0.1)$$

Where $x_{i..}$ indicates the race level of the $ijk$ cell, and $x_{.j.}, x_{..k}$ are defined analogously for the defendant race and disposition. The interaction terms then serve to answer questions about the racial pattern of strikes which is utilized by each party given the defendant race. Most interesting to this investigation is the third order interaction term. This term indicates a significant difference in racial strike patterns given the party striking and the defendant race. In other words, this term accounts for different patterns for the different parties which are not independant of the defendant race.

When this term is tested using a nested model without the third order interaction, the third order interaction is found to be significant. This suggests that not only do the patterns present in the different parties vary, but they vary differently for different defendant races. This dependence can be viewed using a novel graphic presented in Figure 4.1.



Figure 4.1: Parallel coordinate plot of racial strike tendencies

The conditional probability of a particular disposition given the racial combination of venire person and defendant is displayed on the y-axis, that is the count of individuals for a particular race, defendant race, and disposition combination divided by the number of individuals with the racial combination across all dispositions. The x-axis then displays the combinations, grouped by the venire member race to show the dominant pattern in the data.

The black line running across the plot is the mean, or expected, rejection probability that all parties would have if they acted identically. That is, the relative level of this line provides the relative strike rate on aggregate for a particular racial combination. The bars extending from this line at each point go from this line to the corresponding value of the party represented by the bar. Finally, the horizontal lines provide approximate confidence intervals for each combination[1].

The dominant pattern to these strikes is a tendency of the defense to preferentially reject white venire members and keep black venire members, and of the prosecution to do the opposite. It was already noted in the literatureWright et al. (2018), but the addition of defendant race allows us to make a stronger statement, as this pattern remains across defendant races. It also adds nuance, however, as the race of the defendant has a clear impact on the lengths of the bars for both the defense and prosecution. The prosecution seems to favour a jury which does not match the race of the defendant, while the defense seems to favour a jury which does.

While this second tendency seems to have no justification beyond race, the dominant tendency may have other justification than simply skin colour. As was noted by "Ideological Imbalance and Peremptory Challenge", black individuals are more consistently aligned with the democratic party, and as a consequence a lawyer which suspects this political bias will impact the trial outcome would preferentially strike or keep black jurors in order to keep as many left wing individuals as possible. In this data, this political imbalance is incredibly prevalent, as can be seen in Figure 4.2 Add the plot of this effect here, elaborate on this pattern more based on the plot.

Perhaps more interestingly, the prosecution and judge seem to match in their tendency from the mean at every combination. This suggests that both challenges with cause and the prosecution tend to have the same effect on the jury composition, though the magnitudes can differ greatly for these two strikes. An immediate explanation to this is offered by Hans and Vidmar (1986), who outline, on pages 69-70, the skill and tact required to effectively propose challenges with cause. In order to determine an individual's bias, it is frequently the case that a direct question will fail to garner an honest reponse due to social pressures. As a consequence, the questions asked of venire members must be carefully presented.

Using this as a motivation, an obvious possible explanation for the challenges with cause is that the prosecution is simply more experienced on average than the defence. To determine the veracity of this claim, the year licensed for each lawyer was subtracted from the outcome date of each trial. The resulting distribution of years of experience was then plotted in back-to-back histograms as shown in Figure 4.3.

Clearly, this hypothesis is false. It seems the typical defence lawyer is more experienced than the typical prosecutor, not less. Indeed, the prosecutors seem to be much more likely to be inexperienced than the defence lawyers.

---

[1]Generated assuming a binomial distribution of struck (by any party) against kept, as when this data is modelled with a poisson distribution, the distribution of sub-processes given the overall count will be binomially distributed

Figure 4.2: Conditional probabilities of political affiliation by race and gender



Figure 4.3: Distributions of lawyer experience for prosecutors and defence attorneys

# Chapter 5

# Summary

Summarize the presented work. Why is it useful to the research field or institute?

## 5.1 Future Work

Possible ways to extend the work.

# Bibliography

42nd Parliament of Canada (2018a, March). Bill C-75: An Act to Amend the Criminal Code, Youth Criminal Justice Act and other Acts and to make consequential amendments to other Acts. http://www.justice.gc.ca/eng/csj-sjc/pl/charter-charte/c75.html.

42nd Parliament of Canada (2018b, November). Bill C75. LEGISinfo. http://www.parl.ca/LegisInfo.

Attorney General's Office of the United Kingdom (2012, November). Jury vetting right of stand-by guidelines. https://www.gov.uk/guidance/jury-vetting-right-of-stand-by-guidelines–2.

Baldus, D. C., G. Woodworth, D. Zuckerman, and N. A. Weiner (2001). The Use of Peremptory Challenges in Capital Murder Trials: A Legal and Empirical Analysis. *University of Pennsylvania Journal of Consitutional Law 3*(1).

Brown, F. L., F. T. McGuire, and M. S. Winters (1978). The peremptory challenge as a manipulative device in criminal trials: Traditional use or abuse. *New England Law Review 14*, 192.

Brown, R. B. (2000). Challenges for cause, stand-asides, and peremptory challenges in the nineteenth century. *Osgoode Hall Law Journal 38*(3), 453–494.

Burgess, G. (1992). The divine right of kings reconsidered. *The English Historical Review 107*(425), 837–861.

Constitution of Canada (1982). Constitution of Canada. Accessed: https://laws-lois.justice.gc.ca/eng/Const/index.html.

Constitution of the United States (1788). Constitution of the United States. Accessed: https://www.senate.gov/civics/constitution_item/constitution.htm.

Ford, R. (2010). Modeling the effects of peremptory challenges on jury selection and jury verdicts. *George Mason Law Review 17*, 377.

Forsyth, W. (1994). *History of Trial by Jury* (2 ed.). Lawbook Exchange.

Government of Saskatchewan (1998). Jury act, 1998. Accessed: http://www.qp.gov.sk.ca/documents/English/Statutes/Statutes/J4-2.pdf.

Grosso, C. M. and B. O'Brien (2012). A Stubborn Legacy: The Overwhelming Importance of Race in Jury Selection in 173 Post-Batson North Carolina Capital Trials. *Iowa Law Review 97*, 1531.

Hans, V. P. and N. Vidmar (1986). *Judging the Jury* (1 ed.). Plenum Press.

Harris, K. (2018, February). Liberals review jury selection process after Boushie case uproar. CBC News. https://www.cbc.ca/news/politics/jury-selection-diversity-indigenous-1.4531792.

Hasan, N. R. (2018, April). Eliminating peremptory challenges makes trials less fair. The Star. https://www.thestar.com/opinion/contributors/2018/04/10/eliminating-peremptory-challenges-make-trials-less-fair.html.

Hoffman, M. B. (1997). Peremptory Challenges Should Be Abolished: A Trial Judge's Perspective. *The University of Chicago Law Review 64* (3), 809.

Iacobucci, F. (2013). First Nations Representation on Ontario Juries: Report of the Independent Review Conducted by the Honourable Frank Iacobucci. Ministry of the Attorney General. Accessed: https://www.attorneygeneral.jus.gov.on.ca/english/about/pubs/iacobucci/First_Nations_Representatio

Kullback, S. and R. Leibler (1951). On information and sufficiency. *The Annals of Mathematical Statistics 22* (1), 79–86.

MacLean, C. (2018, February). Gerald Stanley acquittal renews calls for justice reform 27 years after Manitoba inquiry. CBC News. https://www.cbc.ca/news/canada/manitoba/aboriginal-justice-inquiry-colten-boushie-gerald-stanley-jury-1.4532394.

Macnab, A. (2018, February). Stanley acquittal should not lead to scrapping peremptory challenges, say criminal lawyers. Canadian Lawyer. https://www.canadianlawyermag.com/legalfeeds/author/aidan-macnab/stanley-acquittal-should-not-lead-to-scrapping-peremptory-challenges-say-criminal-lawyers-15332/.

Ministry of the Attorney General of Ontario (2018). The annual jury selection process. Queen's Printer for Ontario. Accessed: https://www.attorneygeneral.jus.gov.on.ca/english/courts/jury/jury_selection_process.php.

Mirriam-Webster (2019a). Mirriam-Webster Dictionary Online. Accessed: https://www.merriam-webster.com/dictionary/venire.

Mirriam-Webster (2019b). Mirriam-Webster Dictionary Online. Accessed: https://www.merriam-webster.com/dictionary/voir%20dire.

Morehead, J. W. (1994). When a peremptory challenge is no longer peremptory: Batson's unfortunate failure to eradicate invidious discrimination from jury selection. *DePaul Law Review 43*, 625.

North Carolina Sentencing and Policy Advisory Commission (2017). *Classification of a Sample of Offenses.* North Carolina Judicial Branch. Accessed: https://www.nccourts.gov/assets/documents/publications/Sample-list-2017.pdf?MAZluXuS0FWod4nFt7zXecJ8Ifu0qEZx.

Page, A. (2005). Batson's blind spot: Unconscious stereotyping and the peremptory challenge. *Boston University Law Review 85*, 155.

Petersen, C. (1993). Insitutionalized racism: The need for reform of the criminal jury selection process. *McGill Law Journal 38* (1).

Quenneville, G. (2018, February). What happened on Gerald Stanley's farm the day Colten Boushie was shot, as told by witnesses. CBC News. https://www.cbc.ca/news/canada/saskatoon/what-happened-stanley-farm-boushie-shot-witnesses-colten-gerald-1.4520214.

Quenneville, G. and J. Warick (2018, February). Shouts of 'murderer' in courtroom after Gerald Stanley acquitted in Colten Boushie shooting. CBC News. https://www.cbc.ca/news/canada/saskatoon/gerald-stanley-colten-boushie-verdict-1.4526313.

Richardson Oakes, A. and H. Davies (2016). Justice must be seen to be done: a contextual reappraisal. *Adelaide Law Review 37*(2), 461–494.

Roach, K. (2018, April). Ending peremptory challenges in jury selection is a good first step. The Ottawa Citizen. https://ottawacitizen.com/opinion/columnists/roach-ending-peremptory-challenges-in-jury-selection-is-a-good-first-step.

Sheppard, R. (2018). Patriation of the constitution. The Canadian Encyclopedia: Historica Canada. Accessed: https://www.thecanadianencyclopedia.ca/en/article/patriation-of-the-constitution.

Statistics Canada (2018, November). Table 35-10-0061-01: Crime severity index and weighted clearance rates, police services in Saskatchewan.

Supreme Court of Canada (1991). R. v. sherratt. Supreme Court Judgements. SCC Case Number: 21501; Accessed: https://scc-csc.lexum.com/scc-csc/scc-csc/en/item/734/index.do?q=21501.

Supreme Court of the United States (1965). Swain v. Alabama. Accessed: https://supreme.justia.com/cases/federal/us/380/202/.

Supreme Court of the United States (1986). Batson v. Kentucky. Accessed: https://www.law.cornell.edu/supremecourt/text/476/79.

Van Dyke, J. M. (1977). *Jury Selection Procedures: Our Uncertain Commitment to Representative Panels* (1 ed.). Ballinger Publishing.

von Moschzisker, R. (1921). The historic origin of trial by jury. *University of Pennsylvania Law Review 70*(1).

Wickham, H. and J. Bryan (2018). *readxl: Read Excel Files*. R package version 1.1.0.

Woolley, A. (2018). An Ethical Jury? Reflections on the Acquittal of Gerald Stanley for the Murder/Manslaughter of Colten Boushie. *Slaw: Canada's online legal magazine*. Accessed: http://www.slaw.ca/2018/02/20/an-ethical-jury-reflections-on-the-acquittal-of-gerald-stanley-for-the-murder-manslaughter-of-colten-boushie/.

Wright, R. F., K. Chavis, and G. S. Parks (2018, October). The Jury Sunshine Project: Jury Selection Data as a Political Issue. *University of Illinois Law Review 2018*(4), 1407.

Zinchuk, B. (2018, March). Both sides wrong about Stanley trial. Prince George Citizen. https://www.princegeorgecitizen.com/opinion/editorial/both-sides-wrong-about-stanley-trial-1.23199321.

# Appendix A

# Complementary information

Additional material. For example long mathematical derivations could be given in the appendix. Or you could include part of your code that is needed in printed form. You can add several Appendices to your thesis (as you can include several chapters in the main part of your work).

## A.1   Including **R** code with verbatim

A simple (rather too simple, see **??**) way to include code or $R$ output is to use `verbatim`. It just prints the text however it is (including all spaces, "strange" symbols,...) in a slightly different font.

```
## loading packages
library(RBGL)
library(Rgraphviz)
library(boot)

## global variables
X_MAX <- 150

    This allows me to put as many s  p a   c es   as I want.
I can also use \ and ' and & and all the rest that is usually only
accepted in the math mode.

I can also make as
                many
            line
    breaks as
I want... and
            where I want.
```

## A.2  Data Processing Code

However, it is much nicer to use the *listings* package to include R code in your report. It allows you to number the lines, color the comments differently than the code, and so on.

```r
1  ########################################
2
3  ## THESIS DATA PROCESSING SCRIPT
4  ## Christopher Salahub
5  ## Sept 26, 2018
6
7  ########################################
8
9  ## PACKAGES ###########################
10 library(readxl)
11 library(tm)
12 library(stringr)
13 library(grid)
14
15
16 ## CONSTANTS ##########################
17
18 ## start by defining file locations
19 ThesisDir <- "c:/Users/Chris/Documents/ETH Zurich/Thesis/Data"
20 SunshineFile <- paste0(ThesisDir, "/JurySunshineExcel.xlsx")
21 SunshineSheets <- excel_sheets(SunshineFile)
22
23 NorthCarFile <- paste0(ThesisDir,
24                        "/Jury Study Data and Materials/NC Jury Selection Study
                            Database6 Dec 2011.csv")
25
26 PhillyFile <- paste0(ThesisDir,
27                      "/Voir Dire Data & Codebook/capital_venires.csv")
28
29 ## next the factor level codes as given in the codebook and regularized here
30 ## regularization: - political affiliation "N" replaced with "I" for all entries
31 LevRace <- sort(c("A","B","H","N","O","U","W"))
32 LevGen <- sort(c("F","M","U"))
33 LevPol <- sort(c("D","L","R","I","U"))
34
35 ## create a charge tree with regex nodes to identify and clean charge text
36 chargeTree <- list("rape" = list("statutory", "first|1", "second|2"), "sex(?=.*
       offense)" = list("first|1", "second|2"),
37                    "sex(?=.*offend)" = list("regis", "addr"), "murder" = list("
                           first|1" = list("att"), "second|2" = list("att")),
38                    "arson", "firearm" = list("pos", "disch"), "stole" = list("pos
                           "),
39                    "mari" = list("pos", "sell|sale", "man", "pwimsd"), "coca" =
                           list("pos", "sell|sale", "man", "pwimsd"),
40                    "cs" = list("pos", "sell|sale", "man", "pwimsd"), "hero" =
                           list("pos", "sell|sale", "man", "pwimsd"),
41                    "meth" = list("pos", "sell|sale", "man", "pwimsd"),
42                    "oxycod" = list("pos", "sell|sale", "man", "pwimsd"), "mass" =
                            list("pos"), "break" = list("enter"),
43                    "assa" = list("serious bodily", "female", "strangul", "deadly"
                           , "official"),
44                    "larceny" = list("motor", "felon", "merchant"), "false" = list
                           ("pretense"),
45                    "driving" = list("impaired"), "kidnap" = list("first|1", "
                           second|2"),
46                    "robb" = list("dang"), "burg" = list("first|1", "second|2"), "
                           indec" = list("liber"),
47                    "embez", "manslaughter" = list("inv"), "flee" = list("arrest")
                           ,
48                    "abuse|cruelty" = list("child", "anim"), "identity" = list("
                           theft"))
49
50 ## create a list of variables which can sensibly be summarized by trial
```

```
51  TrialVars ← c("TrialNumberID", "DateOutcome", "JudgeID", "DefAttyType", "
        VictimName",
52                  "VictimRace", "VictimGender", "CrimeLocation", "PropertyType",
53                  "ZipCode.Trials", "StateTotalRemoved", "DefenseTotalRemoved",
54                  "CourtTotalRemoved", "JDistrict", "JName", "JRace", "JGender",
55                  "JPoliticalAff", "JVoterRegYr", "JYrApptd", "JResCity", "JResZip",
56                  "ChargeTxt", "Outcome", "Sentence.FullSunshine", "DefendantID.
                      FullSunshine",
57                  "DefendantID.DefendantToTrial", "DefRace", "DefGender", "DefDOB",
                      "DefAttyID",
58                  "DefAttyName", "DCRace", "DCGender", "DCPoliticalAff", "
                      DCYrRegVote",
59                  "DCYrLicensed", "DCResideCity", "DCResideZip", "ProsecutorID", "
                      ProsName",
60                  "ProsRace", "ProsGender", "ProsPoliticalAff", "PYrRegVote", "
                      PYrLicensed",
61                  "PResideCity", "PResideZip", "Guilty", "CrimeType", "DefWhiteBlack
                      ")
62
63
64  ## FUNCTIONS ##########################
65
66  ## Loading and cleaning ###############
67  ## create a descriptive merge function for cleaning (essentially a 'merge'
        wrapper)
68  CleaningMerge ← function(x, y, ...) {
69      ## start by creating the merge
70      ## first match arguments
71      MatchCall ← match.call(merge)
72      MatchCall[[1]] ← quote(merge)
73      ## get input names and ensure proper name structure
74      xname ← MatchCall$x
75      if (!is.symbol(xname)) xname ← as.symbol(paste0(xname[[2]],xname[[3]]))
76      yname ← MatchCall$y
77      if (!is.symbol(yname)) yname ← as.symbol(paste0(yname[[2]],yname[[3]]))
78      ## use this to extract suffixes and fix MatchCall
79      MatchCall$suffixes ← paste0(".", c(xname, yname))
80      MatchCall$x ← xname
81      MatchCall$y ← yname
82      ## specify that the match should be an outer join
83      MatchCall$all ← TRUE
84      ## and use this to make a clean local assignment to modify
85      assign(as.character(xname), cbind(x, Diag.x = 1), envir = environment())
86      assign(as.character(yname), cbind(y, Diag.y = 1), envir = environment())
87      ## now evaluate the call
88      Merged ← eval(MatchCall, envir = environment())
89      ## next perform some checks
90      xExpInds ← is.na(Merged$Diag.x)
91      yExpInds ← is.na(Merged$Diag.y)
92      ## remove the diagnostic columns
93      Merged$Diag.x ← NULL; Merged$Diag.y ← NULL
94      ## summarize the diagnostic checks
95      X_nexp ← sum(xExpInds)
96      Y_nexp ← sum(yExpInds)
97      X_missing ← Merged[xExpInds,]
98      Y_missing ← Merged[yExpInds,]
99      ## print the diagnostics
100     cat("Joined ", paste(xname, yname, sep = " and "), " with ",
101         X_nexp, " and ", Y_nexp, " failed matches respectively \n", sep = "")
102     ## return the results, preferentially keeping the data which is present in x
            but missing from y
103     if (X_nexp == 0 & Y_nexp == 0) {
104         Merged
105     } else list(Merge = Merged[!xExpInds,], Xfails = X_missing, Yfails = Y_
            missing)
106 }
107
108 ## a function to identify and perform swaps with user input
109 SimpleSwapper ← function(data, CorrectLevs, auto = FALSE) {
110     ## first match the data to the columns of interest
```

```
111    colInds ← match(names(CorrectLevs), names(data))
112    ## extract the levels of the columns of interest to check if there are any
           potential swaps
113    swapCheck ← all(sapply(1:length(colInds),
114                           function(ind) identical(sort(levels(as.factor(data[,
                                   colInds[ind]]))),
115                                                   sort(CorrectLevs[[ind]]))))
116    ## if no swaps are present end this check
117    if (swapCheck) {
118        cat("No errors found, exiting.")
119        return(data)
120    }
121    ## if errors are found, further investigate them
122    ## identify potential rows
123    ## first those which have elements out of place
124    SwapPoss ← sapply(1:length(colInds),
125                      function(ind) !(data[,colInds[ind]] %in% CorrectLevs[[ind
                              ]]))
126    ## now rows containing unknown entries
127    Unknown ← sapply(1:length(colInds),
128                     function(ind) data[,colInds[ind]] == "U")
129    ## identify potential swaps by row
130    Swaps ← apply(SwapPoss, 1, function(row) sum(row) > 1)
131    ## identify the potential errors
132    PotErr ← apply(SwapPoss, 1, function(row) sum(row) == 1)
133    ## use the unknowns to account for some errors
134    UnkInd ← apply(Unknown, 1, any)
135    FalErr ← PotErr & UnkInd
136    ## identify the indices to investigate
137    SwapInds ← which(Swaps|FalErr)
138    ErrInds ← which(PotErr & !UnkInd)
139    ## communicate to the user and ask for input
140    cat("There are ", sum(Swaps|FalErr), " swaps to check\n", sep = "")
141    cat("Additionally, it seems there are ", sum(PotErr & !UnkInd), " errors in
           entries\n", sep = "")
142    ## unless automated
143    if (auto) ErrorReturn ← TRUE else ErrorReturn ← as.logical(readline("Return
           the errors? (T/F): "))
144    ## now, if there are possible swaps investigate them
145    if (sum(Swaps|FalErr) != 0) {
146        ## create a temporary storage structure
147        tempRows ← data[SwapInds, colInds]
148        tempRows ← as.data.frame(lapply(tempRows, function(var) levels(var)[as.
               numeric(var)]),
149                                 stringsAsFactors = FALSE)
150        ## loop through and populate this
151        for (ii in 1:nrow(tempRows)) {
152            ## inspect the row
153            print(tempRows[ii,])
154            ## suggest corrections, first generate matches
155            candComb ← lapply(tempRows[ii,],
156                              function(el) which(sapply(CorrectLevs,
157                                                        function(levs) el %in%
                                                                levs)))
158            reps ← unlist(lapply(candComb, length))
159            ## now generate all swap combinations
160            candComb[[1]] ← rep(candComb[[1]], each = max(reps[-1]))
161            candComb ← as.data.frame(candComb, row.names = NULL)
162            ## identify rows which contain all indices, in other words those
                   valid as swaps
163            compRows ← apply(candComb, 1, function(row) all(1:length(CorrectLevs)
                   %in% row))
164            goodComb ← candComb[compRows,]
165            ## clean them up and print them
166            colnames(goodComb) ← NULL
167            rownames(goodComb) ← NULL
168            cat("Potential combinations:\n")
169            print(t(apply(goodComb,1,order)))
170            ## take user input or automatically determine value
171            if (auto) {
```

```
172                         if (!any(compRows)) acceptedComb ← 0 else acceptedComb ← 1
173                     } else acceptedComb ← as.numeric(readline("Enter a combination choice
                           (0 for error, <enter> to accept first): "))
174                     ## handle special cases, 0 if a true error has been identified
175                     if (identical(acceptedComb,0)) { ## 0 if a true error has been
                           identified
176                         ErrInds ← c(ErrInds, SwapInds[ii])
177                         cat("True error identified, adding ", SwapInds[ii], " to error
                               list\n", sep = "")
178                     } else { ## the case where a swap has been correctly identified and
                           selected, or enter has been pressed
179                         ## if enter has been pressed accept the first row
180                         if (is.na(acceptedComb)) acceptedComb ← 1
181                         ## print recombined row
182                         newRows ← tempRows[ii,order(as.matrix(goodComb[acceptedComb,]))]
183                         colnames(newRows) ← NULL
184                         rownames(newRows) ← NULL
185                         cat("Corrected row:")
186                         print(newRows)
187                         cat("------------------\n")
188                         ## correct entry
189                         tempRows[ii,] ← newRows
190                     }
191                 }
192             ## fill the data
193             ## first prevent factor level errors
194             data[,colInds] ← lapply(colInds, function(ind) levels(data[,ind])[as.
                   numeric(data[,ind])])
195             ## now swap the data
196             data[SwapInds,colInds] ← lapply(1:length(colInds), function(ind) tempRows
                   [,ind])
197             ## reconvert back to factors
198             data[,colInds] ← lapply(colInds, function(ind) as.factor(data[,ind]))
199         }
200         ## in either case return the data and errors as specified
201         if (ErrorReturn) {
202             return(list(Data = data, Errors = ErrInds))
203         } else {
204             return(data)
205         }
206 }
207
208 ## now create a function to address the errors possibly identified in the above
        function automatically
209 SwapErrorFix ← function(errorData, CorrectLevs) {
210     ## check if we are in the case without errors
211     if (!identical(names(errorData), c("Data", "Errors"))) {
212         cat("No errors\n")
213         return(errorData)
214     } else {
215         ## extract the data and data in error
216         fulldata ← errorData$Data
217         ## get the relevant columns
218         colInds ← match(names(CorrectLevs), names(fulldata))
219         ## go through the specified variables and remove errors
220         fixed ← lapply(1:length(colInds),
221                     function(ind) {
222                         var ← fulldata[,colInds[ind]]
223                         var ← levels(var)[as.numeric(var)]
224                         inds ← !(var %in% CorrectLevs[[ind]])
225                         cat(names(CorrectLevs)[ind], ": ", sum(inds),
226                             " errors\n", sep = "")
227                         var[inds] ← "U"
228                         as.factor(var)
229                     })
230         ## insert these fixed values
231         fulldata[, colInds] ← fixed
232         ## return this
233         fulldata
234     }
```

```
235  }
236
237  ## write a wrapper to perform this swapping and error correction in one call
238  SwapandError ← function(data, CorrectLevs) {
239      swapped ← SimpleSwapper(data = data, CorrectLevs = CorrectLevs, auto = TRUE)
240      fixed ← SwapErrorFix(errorData = swapped, CorrectLevs = CorrectLevs)
241      fixed
242  }
243
244  ## Variable Synthesis ##################
245  ## Kullback-Leibler divergence function
246  kldiv ← function(samp, dist) {
247      ## convert to matrices
248      mat1 ← as.matrix(samp)
249      mat2 ← as.matrix(dist)
250      ## make into proper distributions
251      mat1 ← mat1/rowSums(mat1)
252      mat2 ← mat2/rowSums(mat2)
253      ## take the log ratio
254      logratio ← log(mat1/mat2)
255      ## multiply by correct matrix
256      vals ← mat1*logratio
257      ## take the row sums
258      rowSums(vals, na.rm = TRUE)
259  }
260
261  ## make a text-mining regularization function
262  StringReg ← function(strs) {
263      ## first set everything to lowercase
264      strs ← tolower(strs)
265      ## replace specific patterns (noticed during early tests)
266      strs ← str_replace_all(strs, "b/e|break/enter|b&e|break or enter|b or e|b &/
             or e|b & e", "breaking and entering")
267      strs ← str_replace_all(strs, "controlled substance", "cs")
268      strs ← str_replace_all(strs, "dwi", "driving while impaired")
269      strs ← str_replace_all(strs, "rwdw", "robbery with a deadly weapon")
270      strs ← str_replace_all(strs, "pwisd|pwmsd|pwmsd|pwitd|pwid|pwmisd|pwosd", "
             pwimsd")
271      strs ← str_replace_all(strs, "robery|rob ", "robbery")
272      strs ← str_replace_all(strs, "bulgary", "burglary")
273      strs ← str_replace_all(strs, "awdw", "assault with a deadly weapon")
274      strs ← str_replace_all(strs, "(?<=[\\sa-z])[0-9]{2,}", "")
275      strs ← str_replace_all(strs, "att ", "attempted ")
276      strs ← str_replace_all(strs, "assult", "assault")
277      strs ← str_replace_all(strs, "marj", "marijuana")
278      ## replace punctuation
279      strs ← gsub("[^[:alnum:][:space:]']", "", strs)
280      ## return these
281      strs
282  }
283
284  ## create a function to process such a tree structure given a list of strings
285  stringTree ← function(strs, regexTree, inds = 1:length(strs), includeOther = TRUE
         ) {
286      ## identify the sublists, and divide the data
287      sublists ← sapply(regexTree, is.list)
288      ## iterate over unnamed items (leaf nodes)
289      listdiv ← lapply(regexTree[!sublists], function(el) inds[grepl(el, strs, perl
             = TRUE)])
290      names(listdiv) ← unlist(regexTree[!sublists])
291      ## check if there are any sublists
292      if (!any(sublists)) {
293          if (includeOther) listdiv ← c(listdiv, other = list(inds[!(inds %in%
                 unlist(listdiv))]))
294          ## in the case of none, treat the object as a list to iterate through
295          listdiv
296      } else {
297          ## otherwise recurse over the branches
298          finlist ← c(listdiv, lapply(names(regexTree)[sublists],
299                                       function(name) stringTree(strs[grepl(name,
```

```
                                                strs, perl = TRUE)],
300                                                           regexTree[[name]],
301                                                           inds[grepl(name,
                                                                  strs, perl =
                                                                  TRUE)],
302                                                           includeOther)))
303           names(finlist)[(length(listdiv) + 1):length(finlist)] ← names(regexTree)[
                  sublists]
304           c(finlist, other = list(inds[!(inds %in% unlist(finlist))]))
305       }
306  }
307
308  ## create a tree depth helper function
309  maxdepth ← function(tree, counter = 1) {
310       max(sapply(tree, function(br) if (!is.list(br)) counter else maxdepth(br,
            counter + 1)))
311  }
312
313  ## create a function to aggregate a tree as specified above at the desired depth
314  treeAgg ← function(tree, level = 1) {
315       ## first check the max depth of the tree
316       treedepth ← maxdepth(tree)
317       ## compare this to requested aggregation level
318       stopifnot(level <= treedepth)
319       ## aggregate at desired level with a helper function
320       agg ← function(tr, depth = 1) {
321           if (depth == level) lapply(tr, function(el) setNames(unlist(el),NULL))
                  else lapply(tr, function(br) agg(dr, depth + 1))
322       }
323       agg(tree)
324  }
325
326  ## create a crime class aggregation function
327  CrimeClassify ← function(tree, regChar) {
328       crimes ← list()
329       crimes$Sex ← unique(c(unlist(tree[c("rape", "sex(?=.*offense)", "sex(?=.*
            offend)", "indec")]),
330                            tree$other[grepl("sex", regChar[tree$other])]))
331       crimes$Theft ← unique(unlist(tree[c("stole", "embez", "break", "larceny", "
            robb", "burg", "identity")]))
332       crimes$Murder ← unique(unlist(tree[c("murder", "manslaughter")]))
333       crimes$Drug ← unique(c(unlist(tree[c("mari", "coca", "cs", "hero", "meth", "
            oxycod")]),
334                            tree$other[grepl("para|drug|substance|pwimsd",
                                 regChar[tree$other])]))
335       crimes$Violent ← unique(unlist(tree[c("arson", "assa", "abuse|cruelty")]))
336       crimes$Driving ← unique(c(unlist(tree[c("driving")]),
337                            tree$other[grepl("hit(?=.*run)|speeding", regChar[
                                 tree$other], perl = TRUE)]))
338       crimes
339  }
340
341  ## in order to make the process of pre-processing the data and adding desired
        columns, place the pre-processing into a
342  ## flexible function and add operations as desired
343  SynCols ← function(data) {
344       ## too busy, synthesize some variables to clearly indicate the results of
            defense and prosecution selection
345       data$VisibleMinor ← data$Race != "White"
346       data$PerempStruck ← grepl("S_rem|D_rem", data$Disposition)
347       data$DefStruck ← data$Disposition == "D_rem"
348       data$ProStruck ← data$Disposition == "S_rem"
349       data$CauseRemoved ← data$Disposition == "C_rem"
350       ## lets look at which race struck each juror
351       data$StruckBy ← as.factor(sapply(1:nrow(data),
352                                          function(ind) {
353                                              dis ← as.character(data$
                                                 Disposition[ind])
354                                              if (dis == "S_rem") {
355                                                  as.character(data$ProsRace
```

```
                                                                            [ind])
356                                                       } else if (dis == "D_rem") {
357                                                            as.character(data$DCRace[
                                                                 ind])
358                                                       } else "Not Struck"
359                                                  }))
360          ## create a white black other indicator
361          data$WhiteBlack ← FactorReduce(data$Race, tokeep = c("Black", "White", "U"))
362          data$DefWhiteBlack ← FactorReduce(data$DefRace, tokeep = c("Black", "White",
                "U"))
363          data$VicWhiteBlack ← FactorReduce(data$VictimRace, tokeep = c("Black", "White
                ", "U"))
364          ## return the data with synthesized columns
365          data
366 }
367
368 ## write functions to process the sentences
369 SentenceProcess ← function(sentencing) {
370          sents ← tolower(sentencing)
371          ## identify sentences in months, years, and days
372          monthsent ← str_extract(sents, "[0-9\\-]+\\s*(?=m)")
373          daysent ← str_extract(sents, "[0-9\\-]+\\s*(?=d)")
374          yearsent ← str_extract(sents, "[0-9\\-]+\\s*(?=y)")
375          ## extract life without parole
376          lwp ← str_extract(sents, "parol[e]*")
377          ## and with parole
378          life ← str_extract(sents, "life")
379          life[!is.na(lwp)] ← NA
380          ## get restitutions
381          resti ← str_extract(sents, "[0-9,]+\\s*(?=restitu)|\\$[0-9,]+")
382          ## get supervised probation
383          suprob ← str_extract(sents, "sup.*pro")
384 }
385
386 ## Summary Functions ###################
387 ## make a function to summarize trial jury data
388 JurySummarize ← function(Varnames = c("Disposition", "Race", "Gender", "
        PoliticalAffiliation")) {
389          ## check if a juror summary object exists already
390          if (!("sun.juror" %in% ls(.GlobalEnv))) {
391               ## first group the data for easy access
392               Juries ← aggregate(sun.swap[, Varnames],
393                                      by = list(TrialNumberID = sun.swap$TrialNumberID,
                                         JurorNumer = sun.swap$JurorNumber),
394                                      unique)
395          } else Juries ← sun.juror
396          ## in either case, perform aggregation by trial instance
397          Juries ← aggregate(Juries[, Varnames],
398                                 by = list(TrialNumberID = Juries$TrialNumberID),
399                                 function(var) var)
400          ## clean up the names
401          names(Juries)[grepl("Polit", names(Juries))] ← "PolAff"
402          Varnames[4] ← "PolAff"
403          ## now summarize relevant features
404          Summary ← apply(Juries[, Varnames], 1,
405                             function(row) {
406                                  ## get final jury indices
407                                  disps ← unlist(row$Disposition)
408                                  foreman ← grepl("Foreman", disps)
409                                  finJur ← grepl("Foreman|Kept", disps)
410                                  defStruck ← grepl("D_rem", disps)
411                                  proStruck ← grepl("S_rem", disps)
412                                  ## process all variables
413                                  newrow ← sapply(row,
414                                                    function(el) {
415                                                         c(Jury = table(unlist(el)[finJur]),
416                                                            Venire = table(unlist(el)),
417                                                            DefRem = table(unlist(el)[
                                                                 defStruck]),
418                                                            ProRem = table(unlist(el)[
```

```
                                                    proStruck]))
419                                    })
420                         newrow$Disposition ← NULL
421                         newrow ← c(unlist(newrow), ForeRace = row$Race[foreman],
422                                   ForeGender = row$Gender[foreman], ForePol =
                                          row$PolAff[foreman])
423                         if (sum(foreman) > 1) {
424                             names(newrow)[names(newrow) == "ForeRace1"] ← "
                                    ForeRace"
425                             names(newrow)[names(newrow) == "ForeGender1"] ← "
                                    ForeGender"
426                             names(newrow)[names(newrow) == "ForePol1"] ← "
                                    ForePol"
427                         }
428                         newrow
429                 })
430     ## perform some clean up
431     longest ← sapply(Summary, length)
432     longest ← which(longest == max(longest))[1]
433     longNames ← names(Summary[[longest]])
434     Summary ← lapply(names(Summary[[longest]]),
435                     function(name) unname(sapply(Summary,
436                                                 function(el) el[name])))
437     names(Summary) ← longNames
438     Summary ← lapply(longNames,
439                     function(nm) {
440                         if (grepl("ForeGender", nm)) {
441                             Summary[[nm]] ← factor(Summary[[nm]], levels = 1:3,
                                    labels = LevGen)
442                         } else if (grepl("ForePol", nm)) {
443                             Summary[[nm]] ← factor(Summary[[nm]], levels = 1:5,
                                    labels = LevPol)
444                         } else if (grepl("ForeRace", nm)) {
445                             Summary[[nm]] ← factor(Summary[[nm]], levels = 1:7,
                                    labels = LevRace)
446                         } else Summary[[nm]]
447                     })
448     names(Summary) ← longNames
449     ## return these
450     list(Juries = Juries, Summaries = as.data.frame(Summary))
451 }
452
453 ## a generic simplification method to summarize a vector
454 Simplifier ← function(col, ...) {
455     UseMethod("Simplifier")
456 }
457
458 ## code up methods for the types to be seen
459 Simplifier.default ← function(col, collapse = "") paste0(col, collapse = collapse
        )
460 Simplifier.numeric ← function(col, na.rm = TRUE, trim = 0, ...) mean.default(col,
        trim = trim, na.rm = na.rm)
461 Simplifier.factor ← function(col, collapse = "", ...) paste0(sort(as.character(
        levels(col)[as.numeric(col)])),
462                                                     collapse = collapse
                                                          )
463 Simplifier.character ← function(col, collapse = "", ...) paste0(sort(col),
        collapse = collapse)
464
465 ## create a grouping wrapper which does unique aggregation of a data set
466 UniqueAgg ← function(data, by, ...) {
467     ## convert data to a data frame for regularity
468     if (!is.data.frame(data)) data ← as.data.frame(data)
469     ## identify the grouping column by in the data
470     by.groups ← names(data) == by
471     ## provide nice error handling
472     stopifnot(sum(by.groups) > 0)
473     ## first identify which rows are already unique
474     groups ← as.numeric(as.factor(unlist(data[by.groups])))
475     unqRows ← sapply(groups, function(el) sum(groups == el) == 1)
```

```
476        ## consider grouping only the other rows using the unique function
477        endata ← data[unqRows,]
478        unqdata ← aggregate(data[!unqRows, !by.groups], by = list(data[!unqRows, by.
               groups]), unique)
479        ## reorder to make sure everything is compatible
480        names(unqdata)[1] ← by
481        unqdata ← unqdata[,match(names(endata), names(unqdata))]
482        ## now use the Simplifier helper defined above to process these results
483        procdata ← lapply(unqdata, function(col) sapply(col, Simplifier, ...))
484        ## append everything together
485        endata ← lapply(1:length(endata),
486                        function(n) c(if (is.factor(endata[[n]])) as.character(
                                endata[[n]]) else endata[[n]],
487                                      procdata[[n]]))
488        names(endata) ← names(data)
489        ## convert to a data frame
490        as.data.frame(endata)
491 }
492
493 ## a simple helper to convert multiple factor levels into a single 'other' level
494 FactorReduce ← function(vals, tokeep) {
495        chars ← as.character(vals)
496        ## simply replace elements
497        chars[!grepl(paste0(tokeep, collapse = "|"), chars)] ← "Other"
498        chars
499 }
500
501 ## write a function to re-level factor variables to make mosaic plots cleaner
502 MatRelevel ← function(data) {
503        temp ← lapply(data, function(el) if (is.factor(el)) as.factor(levels(el)[as.
               numeric(el)]) else el)
504        temp ← as.data.frame(temp)
505        names(temp) ← names(data)
506        temp
507 }
508
509 ## another simple processing function to correct NA's given some other identifier
        and data set
510 FillNAs ← function(dataNAs, filldata, identifier) {
511        ## extract the relevant column indices in a flexible way
512        if (is.null(colnames(filldata))) {
513            relcol ← grepl(identifier, names(filldata))
514        } else relcol ← grepl(identifier, colnames(filldata))
515        ## first identify the relevant rows in the data NAs
516        relRows ← is.na(dataNAs)
517        ## take the relevant rows of the filldata
518        filldata ← matrix(unlist(filldata[relcol]), ncol = sum(relcol))
519        rowfiller ← rowSums(filldata[relRows,])
520        ## return the filled data
521        dataNAs[relRows] ← rowfiller
522        dataNAs
523 }
524
525 ## write a wrapper to estimate the values of total removed jurors
526 RemovedJurorEstimates ← function(tofill, data, ident, plot = TRUE) {
527        temp ← FillNAs(tofill, filldata = data, identifier = ident)
528        temp2 ← rowSums(data[,grepl(ident, names(data))])
529        ## let's see how accurate this is if plotting is desired
530        if (plot) {
531            plot(temp, temp2, xlab = "Observed and Filled", ylab = "Juror Sums")
532            abline(0,1)
533        }
534        cat("= : ", sum(temp == temp2)/length(temp2), "\n", "< : ", sum(temp2 < temp)
               /length(temp2), "\n", sep = "")
535        ## replace the filled values less than the estimated, for consistency
536        temp[temp < temp2] ← temp2[temp < temp2]
537        temp
538 }
539
540
```

```
541 ## LOADING AND PROCESSING DATA #########
542
543 ## load the data
544 SunshineData ← lapply(SunshineSheets, function(nm) as.data.frame(read_excel(
        SunshineFile, sheet = nm)))
545 names(SunshineData) ← SunshineSheets
546 NorthCarData ← read.csv(NorthCarFile)
547 PhillyData ← read.csv(PhillyFile)
548
549 ## clean non-informative columns
550 CleanSunshine ← lapply(SunshineData, function(dat) dat[, !apply(dat,2,function(
        col) all(is.na(col)))])
551
552 ## the Sunshine data needs to be restructured into one table, rather than a
        relational database structure
553 ## see the IDMatch function, this was created specifically to perform ID-based
        table joins
554 ## the most appropriate global target is the juror table, start by matching this
        to the trial
555 FullSunshine ← with(CleanSunshine, CleaningMerge(Jurors, Trials, by = "
        TrialNumberID"))
556 ## remove extra ID column, fix a misleading name
557 FullSunshine$CountyName ← FullSunshine$CountyID
558 FullSunshine$CountyID ← NULL
559 ## clean up two additional columns which had inconsistencies
560 FullSunshine$Disposition ← toupper(FullSunshine$Disposition)
561 FullSunshine$Race[FullSunshine$Race == "?"] ← "U"
562 ## before appending everything to this table, perform some other joins
563 TrialsToCharge ← with(CleanSunshine, CleaningMerge(Charges, Junction, by = "
        ACISID", all = TRUE))
564 DefendantToTrial ← with(CleanSunshine, CleaningMerge(Defendants, DefendantTrial,
        by = "DefendantID", all = TRUE))
565 AttorneyToTrial ← with(CleanSunshine, CleaningMerge(Attorney, AttorneyTrial, by =
         "DefAttyID", all = TRUE))
566 ProsecutorToTrial ← with(CleanSunshine, CleaningMerge(Prosecutor, ProsecutorTrial
        , by = "ProsecutorID", all = TRUE))
567 ## merge issues:
568 ##    - trials to charge: one charge is missing a trial ID, hopefully not
        important
569 ##    - prosecutors to trials: 26 prosecutors without trials, however all entries
         were entirely uninformative
570 ## given the above outputs, rename the failed clean merges to make the next
        section cleaner
571 TrialsToCharge ← TrialsToCharge$Merge
572 ProsecutorToTrial ← ProsecutorToTrial$Merge
573
574 ## now perform some additional merges to create one sheet/data.frame
575 ## add the judge descriptions (no issues)
576 FullSunshine ← CleaningMerge(FullSunshine, CleanSunshine$Judges, by = "JudgeID",
        all = TRUE)
577 ## the charges
578 FullSunshine ← CleaningMerge(FullSunshine, TrialsToCharge, by = "TrialNumberID",
        all = TRUE)
579 ## this leads to 22 jurors in trials without charges and 29 charges without
        trials, inspecting these:
580 ##    - the jurors without charges are all related to a trial with ID number
        "710-01", thankfully the other data
581 ##       for this case is complete, and so it may still be useful for viewing
        jury behaviour
582 ##    - the charges without trials are all of the form "710-0xx", suggesting the
        omission of entire trials of some
583 ##       relation, hopefully these were not too similar, or this exclusion can be
        explained later
584 FullSunshine ← FullSunshine$Merge
585 ## the defendants
586 FullSunshine ← CleaningMerge(FullSunshine, DefendantToTrial, by = "TrialNumberID"
        , all = TRUE)
587 ## the attorneys
588 FullSunshine ← CleaningMerge(FullSunshine, AttorneyToTrial, by = "TrialNumberID",
         all = TRUE)
```

```
589  ## the prosecutors
590  FullSunshine ← CleaningMerge(FullSunshine, ProsecutorToTrial, by = "TrialNumberID
         ", all = TRUE)
591  ## 26 jurors appear to be lacking a prosecutor, these appear to be the
         uninformative prosecutors from earlier, included
592  ## due to the preferential inclusion of the missing values in the first of the
         merged matrices
593  FullSunshine ← FullSunshine$Merge
594
595  ## perform some cleanup
596  ## start with some specific factor replacements
597  ## replace the "N" with "I", as these factor levels are interchangeable in the
         codebook and prevent confusion with race
598  FullSunshine[,grepl("Pol", names(FullSunshine))] ← lapply(FullSunshine[,grepl("
         Pol", names(FullSunshine))],
599                                                        function(var) {
600                                                            var ← toupper(var)
601                                                            var[var == "N"] ←
                                                                    "I"
602                                                            var
603                                                        })
604  ## next save most variables as factors
605  FullSunshine ← lapply(FullSunshine,
606                        function(el) if (is.character(el)) as.factor(el) else el)
607  ## correct some overzealous assignment from above
608  FullSunshine[grepl("Notes", names(FullSunshine))] ← lapply(FullSunshine[grepl("
         Notes", names(FullSunshine))],
609                                                        as.character)
610  ## perform factor regularization according to the factor levels provided in the
         codebook
611  FullSunshine ← sapply(FullSunshine,
612                        function(el) {
613                            if (!is.factor(el)) {
614                                el[el == 999] ← NA
615                                el
616                            } else {
617                                el ← as.character(el)
618                                el ← toupper(el)
619                                el[is.na(el)] ← "U"
620                                as.factor(el)
621                            }
622                        }, simplify = FALSE)
623  FullSunshine ← as.data.frame(FullSunshine)
624  ## remove some unnecessary columns
625  FullSunshine$ID ← NULL
626  FullSunshine$TrialIDAuto ← NULL
627  ## combine the name columns to produce more useful columns
628  FullSunshine$JName ← paste(FullSunshine$JFirstName, FullSunshine$JLastName)
629  FullSunshine$JName[FullSunshine$JName == "U U"] ← "U"
630  FullSunshine$DefAttyName ← paste(FullSunshine$DCFirstName, FullSunshine$
         DCLastName)
631  FullSunshine$DefAttyName[FullSunshine$DefAttyName == "U U"] ← "U"
632  FullSunshine$ProsName ← paste(FullSunshine$ProsecutorFirstName, FullSunshine$
         ProsecutorLastName)
633  FullSunshine$ProsName[FullSunshine$ProsName == "U U"] ← "U"
634
635  ## Checkpoint 1: the clean data has been processed, none of the swaps, synthesis,
          or expansion has taken place
636  ## save this
637  if (!("FullSunshine.csv" %in% list.files())) write.csv(FullSunshine, "
         FullSunshine.csv", row.names = FALSE)
638  ## load if the desire is to start at checkpoint 1
639  if (!("FullSunshine" %in% ls())) FullSunshine ← read.csv("FullSunshine.csv")
640
641  ## Note: the below swap functions have been set to auto as the function's
         performance in these cases has already
642  ## been assessed, and so the swaps have already been inspected, it is critical
         for new data that "auto" be switched
643  ## off to take full advantage of this functionality, and so the wrapper "
         SwapandError" should not be used
```

```
644  ## in the juror data
645  sun.swapJuror ← SwapandError(FullSunshine, CorrectLevs = list(Race = LevRace,
646                                                      Gender = LevGen,
647                                                      PoliticalAffiliation
                                                            = LevPol))
648  ## in the judge data
649  sun.swap ← SimpleSwapper(sun.swapJuror, CorrectLevs = list(JRace = LevRace,
650                                                      JGender =
                                                            LevGen,
651                                                      JPoliticalAff
                                                            = LevPol
                                                            ))
652  ## viewing the error report of these data, they are all related to one judge,
         Arnold O Jones II, who is verified
653  ## as a male after a quick Google search
654  unique(sun.swap$Data[sun.swap$Errors, c("JFirstName", "JLastName")])
655  sun.swapJudge ← sun.swap$Data
656  sun.swapJudge$JGender[sun.swap$Errors] ← "M"
657  sun.swapJudge$JGender ← as.factor(levels(sun.swapJudge$JGender)[as.numeric(sun.
         swapJudge$JGender)])
658  ## in the prosecutor data
659  sun.swap ← SimpleSwapper(sun.swapJudge, CorrectLevs = list(ProsRace = LevRace,
660                                                      ProsGender =
                                                            LevGen,
661                                                      ProsPoliticalAff
                                                            = LevPol
                                                            ))
662  ## that found no errors
663  ## a quick check of the levels of the defendant data finds only one error
664  levels(sun.swap$DefGender)
665  levels(sun.swap$DefRace)
666  sun.swap ← SwapandError(sun.swap, CorrectLevs = list(DefRace = LevRace,
667                                                      DefGender = LevGen
                                                            ))
668  ## next the attorney data
669  sun.swap ← SwapandError(sun.swap, CorrectLevs = list(DCRace = LevRace,
670                                                      DCGender = LevGen,
671                                                      DCPoliticalAff =
                                                            LevPol))
672  ## finally the victim data
673  sun.swap ← SwapandError(sun.swap, CorrectLevs = list(VictimRace = LevRace,
674                                                      VictimGender =
                                                            LevGen))
675  ## this leaves the data error-free (in at least the race/gender/politics columns)
676
677  ## fix the outcome data, which had some improper levels
678  sun.swap$Outcome[sun.swap$Outcome == "HC"] ← "U"
679  sun.swap$Outcome[sun.swap$Outcome == "G"] ← "GC"
680  sun.swap$Outcome ← as.factor(levels(sun.swap$Outcome)[as.numeric(sun.swap$Outcome
         )])
681
682  ## lets make the levels more clear for some of the data (race, politics,
         disposition)
683  ## start with the disposition
684  levels(sun.swap$Disposition) ← c("C_rem", "D_rem", "Foreman", "Kept", "U_rem",
685                                      "S_rem", "Unknown")
686  ## next the political affiliation
687  sun.swap ← lapply(sun.swap, function(el) {
688      if (is.factor(el) & identical(levels(el), LevPol)) {
689          levels(el) ← c("Dem", "Ind", "Lib", "Rep", "U")
690          el
691      } else el})
692  levels(sun.swap$JPoliticalAff) ← c("Dem", "Ind", "Rep", "U")
693  ## now the race
694  sun.swap ← lapply(sun.swap, function(el) {
695      if (is.factor(el) & identical(levels(el), LevRace)) {
696          levels(el) ← c("Asian", "Black", "Hisp", "NatAm", "Other",
697                      "U", "White")
698          el
699      } else el})
```

```
700  levels(sun.swap$VictimRace) ← c("Asian", "Black", "Hisp", "NatAm",
701                                   "U", "White")
702  levels(sun.swap$JRace) ← c("Black", "Hisp", "NatAm", "U", "White")
703  levels(sun.swap$DCRace) ← c("Asian", "Black", "NatAm", "Other",
704                               "U", "White")
705  ## now the outcome/verdict
706  levels(sun.swap$Outcome) ← c("Acquittal", "Guilty as Charged",
707                               "Guilty of Lesser", "Incomplete", "Mistrial",
708                               "U")
709  ## the defense attorney type
710  levels(sun.swap$DefAttyType) ← c("App Priv", "Public", "Private",
711                                   "Ret Priv", "U", "Waived")
712
713  ## add a guilt indicator
714  sun.swap$Guilty ← grepl("Guilty", sun.swap$Outcome)
715
716  ## add a simple indicator of defendant race matching juror race if they are both
         known
717  sun.swap$RaceMatch ← sun.swap$Race == sun.swap$DefRace
718  sun.swap$RaceMatch[sun.swap$Race == "U" | sun.swap$DefRace == "U"] ← NA
719
720  ## now perform tree classification of crimes
721  ## first cast sun.swap as a data frame
722  sun.swap ← as.data.frame(sun.swap)
723  ## regularize the charges
724  chargFact ← as.factor(sun.swap$ChargeTxt)
725  regCharg ← StringReg(levels(chargFact))[as.numeric(chargFact)]
726  ## classify these into a charge tree and aggregate this at the coarsest level
727  aggCharg ← treeAgg(stringTree(regCharg, chargeTree))
728  ## these can be further classified into crime classes
729  crimes.trial ← CrimeClassify(aggCharg, regCharg)
730  ## convert these classes into a factor for the data, start with a generic "other"
          vector
731  sun.swap$CrimeType ← rep("Other", nrow(sun.swap))
732  ## now populate it
733  for (nm in sort(names(crimes.trial))) sun.swap$CrimeType[crimes.trial[[nm]]] ← nm
734  sun.swap$CrimeType ← as.factor(sun.swap$CrimeType)
735
736  ## synthesize additional columns
737  sun.swap ← SynCols(sun.swap)
738
739  ## now organize this on the juror scale
740  sun.juror ← UniqueAgg(sun.swap, by = "JurorNumber", collapse = ",")
741
742  ## Checkpoint 2: the swapped data has been processed and summarized to be on the
         scale of individual jurors
743  ## save the swapped data
744  write.csv(sun.swap, "FullSunshine_Swapped.csv", row.names = FALSE)
745  ## and the juror summarized data
746  saveRDS(sun.juror, "JurorAggregated.Rds")
747
748  ## summarize by trial, get the unique trials
749  Trials ← unique(sun.swap$TrialNumberID)
750  ## extract information about these trials, note that grouping occurs on the trial
         ID, defendant ID, and charge ID levels,
751  ## as the trials frequency involve multiple charges and defendants, which makes
         them less clean
752  sun.trial ← aggregate(sun.swap[,TrialVars],
753                        by = list(sun.swap$TrialNumberID, sun.swap$DefendantID
                                 .DefendantToTrial,
754                                 sun.swap$ID.Charges),
755                        unique)
756  sun.trial$Group.1 ← NULL
757  sun.trial$Group.2 ← NULL
758  sun.trial$Group.3 ← NULL
759
760  ## summarize the juries by trial as well
761  sun.jursum ← JurySummarize()
762
763  ## merge the summaries to the trial sunshine data
```

```
764  sun.trialsum ← merge(cbind(TrialNumberID = sun.jursum$Juries$TrialNumberID, sun.
         jursum$Summaries),
765                          sun.trial, all = TRUE)
766
767  ## notice that the total removed variables are incomplete, try to correct this
         where possible using the jury
768  ## summarized data above
769  sun.trialsum$DefRemEst ← RemovedJurorEstimates(sun.trialsum$DefenseTotalRemoved,
         data = sun.trialsum,
770                                                  ident = "Gender.DefRem", plot =
                                                      FALSE)
771  ## perform this same procedure for the prosecution removals
772  sun.trialsum$ProRemEst ← RemovedJurorEstimates(sun.trialsum$StateTotalRemoved,
         data = sun.trialsum,
773                                                  ident = "Gender.ProRem", plot =
                                                      FALSE)
774  ## synthesize some other variables, simple race indicators
775  sun.trialsum$DefWhiteBlack ← as.factor(FactorReduce(sun.trialsum$DefRace, tokeep
         = c("Black", "White", "U")))
776  sun.trialsum$DefWhiteOther ← as.factor(FactorReduce(sun.trialsum$DefWhiteBlack,
         tokeep = c("White", "U")))
777  ## the Kullback-Leibler divergence
778  sun.trialsum$KLdiv ← kldiv(sun.trialsum[,grepl("Jury", names(sun.trialsum))],
779                             sun.trialsum[,grepl("Venire", names(sun.trialsum))])
780
781  ## Checkpoint 3: the data has been set to the trial level and summarized
782  ## save this
783  saveRDS(sun.trialsum, "TrialAggregated.Rds")
784  saveRDS(sun.jursum, "AllJuries.Rds")
785
786
787  ## CHARGE CLASSIFICATION IMAGES ########
788  ## presented here is the code to generate the appendix charge classification
         images
789
790  ## extract relevant charges from the clean sunshine data, avoiding duplicates
791  ## regularize the charges
792  chargFact.clean ← as.factor(CleanSunshine$Charges$ChargeTxt)
793  regCharg.clean ← StringReg(levels(chargFact.clean))[as.numeric(chargFact.clean)]
794  ## classify these into a charge tree and aggregate this at the coarsest level
795  tree.clean ← stringTree(regCharg.clean, chargeTree)
796
797  ## define padding
798  crimepad ← 0.01
799  areatop ← 0.94
800  hght ← ((0.94 - 0.035) - 5*crimepad)/5
801  wid ← ((0.975 - 0.025) - 7*crimepad)/6
802
803  ## add overall box
804  grid.newpage()
805  grid.rect(width = 0.95, height = 0.95)
806  grid.text(label = "Charges", x = 0.025 + crimepad/2, y = 0.975 - crimepad/2, just
         = c("left","top"))
807  grid.text(label = "(198)", x = 0.975 - crimepad/2, y = 0.975 - crimepad/2, just =
         c("right","top"))
808
809  ## inner crime boxes
810  xpos ← crimepad + 0.025 + wid/2 + (0:5)*(crimepad + wid)
811  ypos ← crimepad + 0.025 + hght/2 + (0:4)*(crimepad + hght)
812  coords ← cbind(rep(xpos, each = 5), rev(ypos))[1:length(tree.clean)-1,]
813  grid.rect(x = coords[,1], y = coords[,2], width = wid, height = hght)
814  grid.text(label = names(tree.clean)[1:(length(tree.clean)-1)], x = coords[,1] - (
         wid - crimepad)/2,
815          y = coords[,2] + (hght-crimepad)/2, just = c("left","top"))
816
817  ## add unclassified and top level counts
818  unclass ← sapply(tree.clean, function(el) if (is.list(el)) length(el$other) else
         0)
819  grid.text(label = paste0("(", unclass, ")"), x = coords[,1] + (wid - crimepad)/2,
         y = coords[,2] + (hght-crimepad)/2,
```

```
820              just = c("right", "top"))
821 toplev ← sapply(tree.clean, function(el) if (is.list(el)) "" else length(el))
822 grid.text(label = toplev, x = coords[,1], y = coords[,2])
823
824 ## add sublist counts
825 for (ii in 1:length(tree.clean)) {
826     if (is.list(tree.clean[[ii]])) {
827         currlist ← tree.clean[[ii]]
828         len ← length(currlist) - 1
829         boty ← coords[ii,2] - hght/2
830         newwid ← wid - 2*crimepad
831         newhght ← (hght - 0.035 - len*crimepad)/len
832         newy ← crimepad + newhght/2 + (0:(len-1))*(newhght + crimepad) + boty
833         if (names(tree.clean)[ii] == "assa") {
834             grid.rect(x = coords[ii,1], y = newy, width = newwid, height =
                     newhght)
835             grid.text(label = names(currlist), x = coords[ii,1] - (newwid-
                     crimepad)/2, y = newy,
836                     just = "left", gp = gpar(fontsize = 8))
837             grid.text(sapply(currlist, length), x = coords[ii,1], y = newy, gp =
                     gpar(fontsize = 8))
838         } else if (names(tree.clean)[ii] == "murder") {
839             grid.rect(x = coords[ii,1], y = newy, width = newwid, height =
                     newhght)
840             grid.text(label = names(currlist), x = coords[ii,1] - (newwid-
                     crimepad)/2, y = newy + (newhght-crimepad)/2,
841                     just = c("left","top"), gp = gpar(fontsize = 8))
842             grid.text(paste0("(", sapply(currlist[1:len], function(el) length(el$
                     other)), ")"), x = coords[ii,1] + (newwid-crimepad)/2,
843                     y = newy + (newhght-crimepad)/2, just = c("right","top"),
                         gp = gpar(fontsize = 8))
844             grid.rect(x = coords[ii,1], y = newy - 0.01, width = newwid - 2*
                     crimepad, height = (newhght - 0.02)/2)
845             grid.text(rep("att",2), x = coords[ii,1] - (newwid-2*crimepad-
                     crimepad)/2, y = newy-0.01,
846                     just = "left", gp = gpar(fontsize = 8))
847             grid.text(sapply(currlist[1:len], function(el) length(el$att)), x =
                     coords[ii,1],
848                     y = newy - 0.01, gp = gpar(fontsize = 8))
849         } else {
850             grid.rect(x = coords[ii,1], y = newy, width = newwid, height =
                     newhght)
851             grid.text(label = names(currlist), x = coords[ii,1] - (newwid-
                     crimepad)/2, y = newy + (newhght-crimepad)/2,
852                     just = c("left","top"), gp = gpar(fontsize = 8))
853             grid.text(sapply(currlist, length), x = coords[ii,1], y = newy, gp =
                     gpar(fontsize = 10))
854         }
855     }
856 }
857
858 ## a small example tree
859 firstx ← 0.33
860 secondx ← 0.75
861 firsty ← c(0.25,0.75)
862 secondy ← c(0.125,0.375,0.625,0.875)
863 wd ← 0.3
864 hg ← 0.1
865 grid.newpage()
866 grid.rect(x = firstx, y = firsty, width = wd, height = hg)
867 grid.rect(x = secondx, y = secondy, width = wd, height = hg)
868 grid.lines(x = c(0,firstx-wd/2), y = c(0.5,firsty[1]))
869 grid.lines(x = c(0,firstx-wd/2), y = c(0.5,firsty[2]))
870 grid.lines(x = c(firstx+wd/2,secondx-wd/2), y = c(firsty[2],secondy[3]))
871 grid.lines(x = c(firstx+wd/2,secondx-wd/2), y = c(firsty[2],secondy[4]))
872 grid.lines(x = c(firstx+wd/2,secondx-wd/2), y = c(firsty[1],secondy[1]))
873 grid.lines(x = c(firstx+wd/2,secondx-wd/2), y = c(firsty[1],secondy[2]))
874 grid.text(label = c("sex(?=.*offense)", "sex(?=.offend)"), x = firstx, y = firsty
        )
875 grid.text(label = c("first|1","second|2","regis","addr"), x = secondx, y = rev(
```

```
    secondy))
```

## A.3   Jury Sunshine Irregularities

Table A.1: Jury sunshine data irregularities noted in data flattening

| Charges without trial (ACISID) | 08CRS50940, 08CRS52888, 09CRS000305, 09CRS1106, 09CRS50752, 10CR52031, 10CRS051975, 10CRS1215, 10CRS397, 10CRS51388, 10CRS51610, 10CRS52410, 11CRS051642, 11CRS051795, 11CRS1577, 11CRS1745, 11CRS1783, 11CRS51204, 11CRS51895, 11CRS52470, 08CRS54836, 08CRS50113 |
|---|---|
| Prosecutors without trials (IDs) | 1-000, 11B-000, 12-000, 14-000, 15B-000, 16A-000, 16B-000, 17A-000, 17B-000, 19A-000, 19B-000, 20A-000, 20B-000, 21-000, 22A-000, 22B-000, 24-000, 25-000, 27A-000, 27B-000, 28-000, 29A-000, 29B-000, 30-000, 6-000, 9-000 |
| Trial missing charge (ID) | 710-01 |

## A.4   Jury Sunshine Charge Classification

## A.5   Using Sweave to include **R** code (and more) in your report

The easiest (and most elegant) way to include R code and its output (and have all your figures up to date with your report) is to use Sweave. You can find an introduction Sweave in `/u/sfs/StatSoftDoc/Sweave/Sweave-tutorial.pdf`.

**Charges (198)**

- **arson (0):** 5
- **murder (8):** second|2 (14): att 2; first|1 (55): att 17
- **cs (14):** pwimsd 12; man 3; sell|sale 8; pos 10
- **break (1):** enter 136
- **kidnap (4):** second|2 12; first|1 16
- **flee (1):** arrest 18
- **embez (0):** 17
- **firearm (9):** disch 5; pos 50
- **hero (5):** pwimsd 1; man 0; sell|sale 1; pos 3
- **assa (31):** official 8; deadly 64; strangul 9; female 11; serious bodily 18
- **robb (41):** dang 58
- **abuse|cruelty (0):** anim 3; child 11
- **rape (4):** second|2 18; first|1 22; statutory 12
- **stole (4):** pos 36
- **meth (3):** pwimsd 0; man 3; sell|sale 0; pos 6
- **larceny (41):** merchant 5; felon 26; motor 11
- **burg (1):** second|2 9; first|1 18
- **identity (0):** theft 7
- **sex(?=.*offense) (5):** second|2 7; first|1 30
- **mari (5):** pwimsd 23; man 4; sell|sale 13; pos 12
- **oxycod (1):** pwimsd 1; man 0; sell|sale 0; pos 0
- **false (5):** pretense 39
- **indec (5):** liber 28
- **sex(?=.*offend) (4):** addr 9; regis 15
- **coca (20):** pwimsd 44; man 1; sell|sale 30; pos 51
- **mass (0):** pos 4
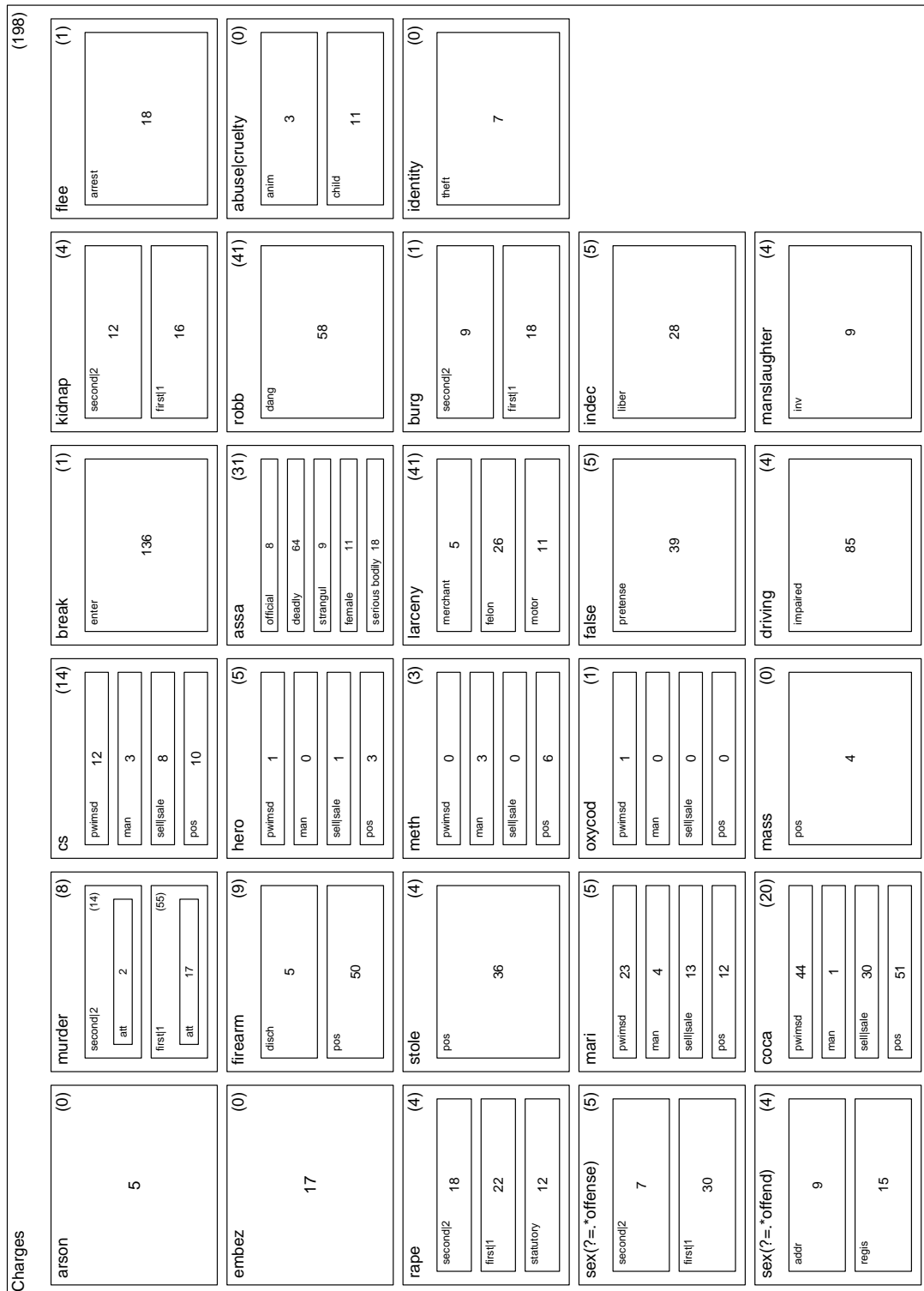- **driving (4):** impaired 85
- **manslaughter (4):** inv 9

Figure A.1: The regular expression charge tree arranged by hierarchy with counts provided. The counts in brackets indicate the counts of charges which could not be classified to a lower level of the hierarchy

# Appendix B

# Yet another appendix....

## B.1 Description

**Something** details.

**Something else** other definition.

## B.2 Tables

Refer to Table B.1 to see a left justified table with caption on top.

Table B.1: Results.

| Student | Grade |
|---------|-------|
| Marie | 6 |
| Alain | 5.5 |
| Josette | 4.5 |
| Pierre | 5 |

# Epilogue

A few final words.

# Declaration of Originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

---

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor .

**Title of work** (in block letters):

> . . .

**Authored by** (in block letters):
*For papers written by groups the names of all authors are required.*

| **Name(s):** | **First name(s):** |
| --- | --- |
| *Mustern* | *Student* |
| | |
| | |
| | |
| | |

With my signature I confirm that
- I have committed none of the forms of plagiarism described in the Citation etiquette information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work .
- I am aware that the work may be screened electronically for plagiarism.
- I have understood and followed the guidelines in the document *Scientific Works in Mathematics*.

**Place, date:**                              **Signature(s):**

*Zurich August 19th 2009*                     *bla*

*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*