

Submitted to *Manufacturing & Service Operations Management*
manuscript (Please, provide the manuscript number!)

Authors are encouraged to submit new papers to INFORMS journals by means of a style file template, which includes the journal title. However, use of a template does not certify that the paper has been accepted for publication in the named journal. INFORMS journal templates are for the exclusive purpose of submitting to an INFORMS journal and should not be used to distribute the papers in print or online or to submit the papers to another publication.

Analyzing the Effects of Judicial Rotation on Criminal Sentencing: An Operations Management Perspective

Baris Ata

Booth School of Business, University of Chicago, Baris.Ata@chicagobooth.edu

Rhys Hester

College of Behavioral, Social and Health Sciences, Clemson University, rhysh@clemson.edu

Lawrence M. Wein

Graduate School of Business, Stanford University, lwein@stanford.edu

Yuwei Zhou

Booth School of Business, University of Chicago, Yuwei.Zhou@chicagobooth.edu

Problem definition: To understand how the impact of judicial rotation and subsequent judge shopping on the defendant's sentence length is mediated by three operational characteristics: the amount of judicial rotation, the allowable shopping time window for defendants, and the capacity utilization of the judicial system.

Methodology/results: Using data from South Carolina in 2000-2001, we formulate and calibrate a mathematical model in which judges rotate across counties, defendants shop for judges, and the sentencing (either by plea or trial) is the result of strategic interactions among the defendant, the judge and the prosecutor. We vary the three operational characteristics via simulation.

The mean and standard deviation of the defendant sentence length decreases (with decreasing returns to scale) in the amount of judicial rotation and the allowable shopping window for defendants, and increases in the capacity utilization, with judicial rotation and the shopping window exhibiting synergistic behavior. The average reduction is modest ($\leq 10\%$), although a small proportion of defendants are impacted in a significant way. In a variant of the model adapted to an urban setting where all defendants have access to all judges, the mean and standard deviation of the sentence length decreases in the number of judges, even in the absence of intertemporal judge shopping.

Managerial implications: Judicial rotation in a rural setting can lead to a modest reduction in the mean sentence and to more equitable sentencing. These effects can occur naturally in an urban setting.

Key words: Plea bargaining, inequality, queuing

1. Introduction

This study is motivated by the empirical findings in Hester (2017), which looks at criminal sentencing outcomes in South Carolina during 2000-2001. A distinctive aspect of South Carolina's judicial system, which consisted of 50 judges presiding over 46 counties, was judicial rotation: although judges spent much of their time in their home circuit, they traveled to an average of 12 counties (typically holding court for a week in each county) and counties encountered an average of 13 different judges throughout the year. As is the case in other U.S. jurisdictions, more than 98% of sentenced cases ended in a plea bargain and less than 2% went to trial. Results in Hester (2017) revealed that judicial rotation had two behavioral effects: it led to judge shopping, where defendants would strategically wait for a lenient judge before agreeing to a plea bargain, and to the cross-pollination of ideas and norms by increasing the interactions among judges and prosecutors. Consequently, judicial rotation led to a decrease in both the mean and standard deviation of the sentence length in South Carolina. Indeed, even though South Carolina was a non-guidelines state, the inter-county variability in sentence lengths was smaller than in most guideline states. See Hester (2017) for a fuller discussion of the contextual setting and the results.

In this study, we construct a mathematical model that attempts to capture the effects of judge shopping but does not explicitly incorporate cross-pollination among judges. We model the process in which defendants shop for a judge, the prosecutor proposes a plea deal to the defendant, where both sides are aware of the leniency of the chosen judge, and the defendant either accepts the plea deal or goes to trial. A distinctive aspect of our model is the consideration of queueing and congestion: arriving defendants can only choose judges that have slack capacity in their schedule.

There were not sufficient data to perform an econometric analysis of the interactions among the judge, prosecutor and defendant. Specifically, while we know which judge was chosen by each defendant, we do not know which judges were previously rejected by the defendant during the shopping process. Nonetheless, we are able to estimate the model parameters from the South Carolina data. By simulating the mathematical model, we compute the mean and standard deviation of the sentence length as a function of three operational characteristics: the amount of judicial rotation, the allowable shopping time window for defendants, and the capacity utilization of the judicial system (i.e., the total case workload divided by the total judicial capacity).

Because the full statewide judicial rotation scheme of South Carolina appears unique to that jurisdiction, we also adapt the model to an urban setting where defendants have the opportunity to shop for judges. This context of multiple judges within an urban jurisdiction has been the more typical application of the ideas of individual versus master calendaring, and likely has more far-reaching implications in larger jurisdictions where parties can attempt strategically to identify favorable judges. Applying our model in this setting helps to quantify the effects of shopping availability and capacity utilization in large jurisdictions.

2. Literature Review

By some accounts, the identity of the sentencing judge may matter more to a case outcome than the facts of the case or the background of the defendant. A robust line of sentencing literature focuses on the importance of the identity of the sentencing judge (e.g., Frazier and Bock (1982); Johnson (2006); Myers and Talarico (1987); Spohn (1990); Steffensmeier and Britt (2001)). Consistent with Ulmer (2019) application of Inhabited Institutions Theory to the study of courts and sentencing, we consider the broader impact that court infrastructure characteristics can have on outcomes through the mechanism of the judicial calendaring system. Inhabited Institutions Theory emphasizes how individual actors exercise discretion, reacting to and contributing institutional rules and cultures. In our context, we consider how actors are able to use judge assignment rules to effect more optimal sentencing outcomes—that is, how parties are able to use calendaring systems and local rules and norms to strategically shop for more favorable judges. Early work by Eisenstein et al. (1988) and Ulmer (1997) highlighted differences between individual calendaring systems and master calendaring systems. Under individual calendaring systems judges are assigned a case and retain control of it while in master calendaring systems different judges may handle various tasks, such as arraignment, motions, presiding over the plea or trial, and sentencing. Because the style of calendaring system can be a matter of local rules, systems may differ across counties even within the same jurisdiction (Ulmer (1997)). This early work on court communities found that in master calendaring systems, parties were sometimes able to influence court administrators or otherwise strategically engage in motions or delays in order to “judge shop.” Hester (2017: 218) found that South Carolina’s statewide master calendaring system in conjunction with the practice of judicial rotation led to “an exaggerated form of shopping” for judges. Using a mixed methods

approach he found the practice of regular judge rotation (in which judges routinely traveled from county to county holding court) led to the influence of “plea judges”—lenient judges whose sentencing preferences established baseline norms or going rates for sentencing. Since judges rotated, defendants could strategically choose to enter guilty pleas when plea judges were holding court in their jurisdiction. This reality led other pragmatic judges to adopt plea judge norms for the sake of efficiency. We extend Hester (2017) by formulating and calibrating a mathematical model of judicial rotation. We also adapt this model to settings involving defendants in a large urban setting, which may offer insights into how our formal model of rotation in South Carolina would generalize to other settings involving judge shopping.

The bulk of our modeling involves the dynamic judge shopping process, which is influenced by the recent models in Yang (2016), Silveira (2017) and Wang (2019). These are part of a much larger literature that analyzes models of the plea versus trial game between a defendant and a prosecutor (e.g., Gross and Syverud (1991) and references therein). Among the three operational characteristics we study, only capacity constraints have been modeled in the judicial context (e.g., Ostrom et al. (1999)).

3. Model

Our simulation model of the judicial process involves three agents for each case – the judge, the prosecutor and the defendant – and results in either a plea bargain or a trial. The output of the simulation model includes three primary performance measures that are defined in §3.4: the mean and the standard deviation of the sentence length of all plea bargains (the mean sentence at trials is assumed to be the same for each judge, and hence is omitted) and the standard deviation of the mean sentence length across counties. The goal of our analysis is to understand how these performance measures are impacted by three key operational characteristics: the amount of judicial rotation, the defendant shopping window, and the aggregate judge utilization.

Our simulation model consists of three modules that are described in §3.1-3.3: the defendant arrivals, the judge schedule and the plea bargain, with each module containing one of the three key operational characteristics. The defendant arrivals module describes the timing and characteristics of defendants that arrive to the system, where the arrival rate depends upon the specified judge utilization. The judge schedule module assigns judges to

counties each week based on the specified amount of judicial rotation, and computes the number of pleas that they can process. The plea bargain module includes the interactions among the three agents, and the final sentence imposed on a defendant is determined by the judge he chooses (which in turn depends on the defendant shopping window), the plea offer recommended or presented by the prosecutor and approved (or specified, if a straight plea) by the judge, and whether the defendant accepts the plea offer (i.e., an agreement is reached) or refuses the plea offer and goes to trial. Figure 1 shows a high-level description of our model structure. The estimation of the parameters for each module is described in §4. For ease of reference, all model parameters are described in Table 1.

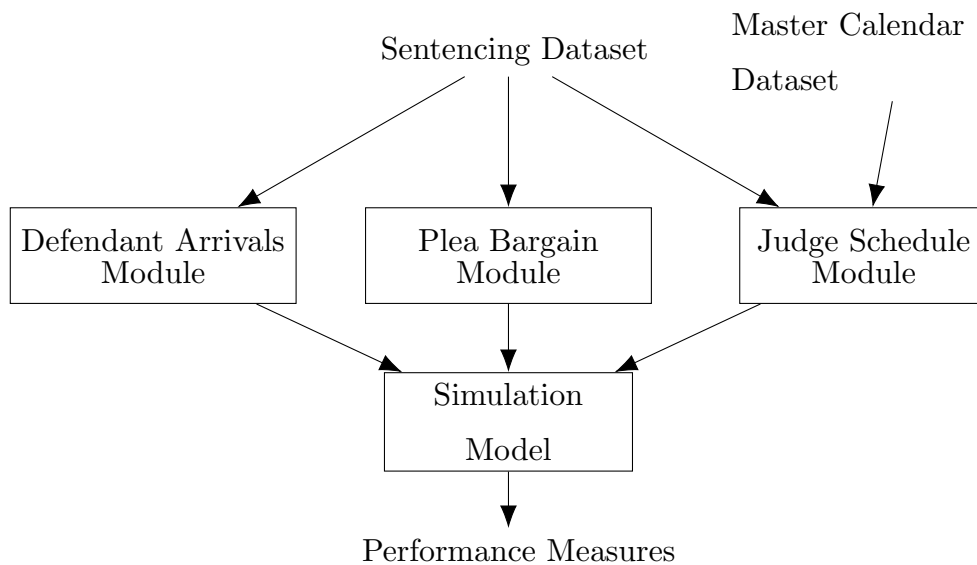


Figure 1 The simulation model contains three modules: the defendant arrivals (§3.1), the judge schedule (§3.2) and the plea bargain (§3.3). The performance measures are defined in §3.4, and the two datasets are described in §4.1.

3.1. Defendant Arrivals

Defendants in our model are indexed by $i = 1, 2, \dots$ and each arriving defendant is assigned a set of covariate values \mathbf{x}_i , which for ease of presentation is suppressed and embedded in the subscript i . The number of defendants who arrive to county c during a week is modeled by a Poisson random variable N_c with mean λ_c . The value of λ_c is dictated by the key operational characteristic ρ , which is the desired judge utilization (i.e., the overall proportion of available time that judges spend presiding over cases) in our simulation

Table 1 The model parameters

Parameter	Description	Module	Value
ρ	Judge utilization	Defendant arrivals	Specified
λ_c	Defendant arrival rate		(6)
\mathbf{x}_i	Defendant covariates		§4.3, Table 2
η	Judge travel probability	Judge schedule	Specified
T_{jct}	Available hours for judge j to work in county c in week t		Master calendar
n_{jct}^T	Number of trials in county c assigned to judge j in week t		Master calendar
γ_P	Mean processing time for a plea		0.018 weeks
γ_T	Mean processing time for a trial		0.839 weeks
r	Defendant shopping window	Plea bargain	Specified
p_i	Probability defendant i receives a zero-length plea sentence		(7), Table EC.3
θ_i	Probability of conviction at trial for defendant i		(8), Table EC.4
τ_i	Expected sentence if convicted at trial for defendant i		(9), Table EC.5
$\ell_j(\cdot)$	Lower bound on the approved sentence length from judge j		§4.4, Fig. EC.4
$u_j(\cdot)$	Upper bound on the approved sentence length from judge j		§4.4, Fig. EC.4
$c_d(i)$	Trial cost for defendant i (in months)		§4.4, Fig. EC.6
d	Defendant waiting cost per week (in months per week)		0.1

model. In §4.3, we compute λ_c , see Equation (6), and describe the assignment of the covariate values \mathbf{x}_i .

3.2. Judge Schedule

We index judges by $j = 1, \dots, J$ and counties by $c = 1, \dots, C$. The judge schedule module specifies the location (county or counties) of judge j in each week t , along with the number of plea cases that judge j can process in week t . Each judge is assigned one (and occasionally more than one) home county, and the weekly judge locations are driven by the key operational characteristic η , which is the probability that a judge is traveling (i.e., working at a non-home county) in any given week. The weekly judge locations in the module output are denoted by C_{jct} , which is the proportion of judge j 's time in week t that is assigned to county c .

For $j = 1, \dots, J$ and $t = 1, \dots, 52$ in the master calendar dataset, let T_{jct} be the proportion of week t that judge j is available to process pleas and trials in county c , and let n_{jct}^T be the number of trials in county c assigned to judge j in week t . Let γ_P and γ_T be the mean processing time (independent of judge) for a plea and a trial, respectively. Then the number of pleas that judge j can process in week t in county c is denoted by N_{jct} and modeled as a Poisson random variable with mean $(T_{jct} - n_{jct}^T \gamma_T)^+ / \gamma_P$.

The parameters T_{jct} , n_{jct}^T , γ_P and γ_T are estimated in §4.2, and η is the specified operational characteristic. The outputs of the judge schedule module are C_{jct} and N_{jct} .

3.3. Plea Bargain Process

The key operational characteristic in the plea bargain module is the defendant shopping window r , which is an integer number of weeks. An arriving defendant chooses the available judge within his shopping window that minimizes his total cost. In our model, defendant i receives a zero-length sentence (e.g., probation or supervision) with probability p_i , which is a function of the defendant's covariates but (see §4.4 for a justification) independent of the presiding judge.

We now describe the plea bargain process under the assumption that defendant i does not receive a zero-length sentence. Defendants may reject a plea deal and decide to go to trial. Depending on a defendant's covariates \mathbf{x}_i , if he goes to trial then he is convicted with probability θ_i , in which case he receives the expected sentence length τ_i . In addition, if the defendant goes to trial, which is a longer judicial process, he incurs the additional cost $c_d(i)$, which is in time units. Hence, the expected total cost of resolving defendant i 's case via trial is $\theta_i\tau_i + c_d(i)$, which is observable by all three agents.

Crucially, judges vary in their leniency, which is modeled using two leniency functions $\ell_j(\cdot)$ and $u_j(\cdot)$ for each judge j . These functions define a lower and upper bound, respectively, on the sentence length as a function of $\theta_i\tau_i$, which is the unconditional mean sentence length if defendant i goes to trial. More specifically, judge j only approves plea deals with a sentence length in the range $[\ell_j(\theta_i\tau_i), u_j(\theta_i\tau_i)]$, where $\ell_j(\theta_i\tau_i) \leq u_j(\theta_i\tau_i)$ for all j and for all values of $\theta_i\tau_i \geq 0$.

We assume that prosecutors try to avoid trials by resolving cases through plea bargaining. However, in doing so, prosecutors aim to maximize the sentence length, leading to three possible outcomes:

- If $\theta_i\tau_i + c_d(i) > u_j(\theta_i\tau_i)$, then the prosecutor offers $u_j(\theta_i\tau_i)$, which the defendant accepts, and the judge approves the plea.
- If $\theta_i\tau_i + c_d(i) < \ell_j(\theta_i\tau_i)$, then there are no plea offers that both the defendant would accept and the judge would approve. The case goes to trial and the defendant's expected cost is $\theta_i\tau_i + c_d(i)$.
- If $\ell_j(\theta_i\tau_i) \leq \theta_i\tau_i + c_d(i) \leq u_j(\theta_i\tau_i)$, then the prosecutor offers $\theta_i\tau_i + c_d(i)$, which the defendant accepts and the judge approves.

To summarize, $\min\{\theta_i\tau_i + c_d(i), u_j(\theta_i\tau_i)\}$ is the expected sentence that defendant i receives either through a plea bargain or at trial, provided that he is not offered a zero-length

sentence during the plea bargaining process, which has probability p_i . To pick judge j , the defendant needs to delay going to court until judge j works in the defendant's county and has sufficient capacity to hear the case. Let d denote the defendant's cost of waiting per week and $w_i(j)$ denote the number of weeks that he needs to wait for judge j . Recall that defendants are given the window of r weeks to shop for judges. Letting $J_i(r)$ denote the set of available judges, i.e., judges who have assigned sessions with remaining capacity within r weeks upon defendant i 's arrival, the defendant chooses judge j^* such that

$$j^* \in \arg \min_{j \in J_i(r)} [(1 - p_i) \min \{\theta_i \tau_i + c_d(i), u_j(\theta_i \tau_i)\} + w_i(j)d]. \quad (1)$$

The parameters p_i , θ_i , τ_i , $c_d(i)$, $\ell_j(x)$, $u_j(x)$ and d are estimated in §4.4, whereas the shopping window r is a key operational characteristic that we specify. The parameter $w_i(j)$ and the set $J_i(r)$ are dictated by the weekly judge-location assignments C_{jct} and the weekly judge capacities N_{jct} , both of which are outputs of the judge schedule module. That is, we make judge j unavailable in week t after she has been assigned N_{jct} plea cases, and update the $J_i(r)$ sets accordingly. The output of the plea bargain module is the expected (nonzero) sentence for each defendant i , $\min \{\theta_i \tau_i + c_d(i), u_{j^*}(\theta_i \tau_i)\}$.

3.4. Performance Measures

Because the probability of a zero-length sentence is independent of the judge in our model, our three performance measures include only cases with nonzero sentences. Let I_c be the set of simulated defendants from county c whose sentence length is positive, and let $I = \cup_{c=1}^C I_c$. Let $s_i = \min \{\theta_i \tau_i + c_d(i), u_{j^*}(\theta_i \tau_i)\}$ be defendant i 's sentence length in the model output. The mean and standard deviation of the plea sentence length are given by

$$\mu = \frac{\sum_{i \in I} s_i}{|I|} \text{ and } \sigma = \sqrt{\frac{\sum_{i \in I} (s_i - \mu)^2}{|I|}}.$$

Our final performance measure is the standard deviation across counties of the mean plea sentence length:

$$\sigma_c = \sqrt{\frac{\sum_{c=1}^C (\mu_c - \mu)^2}{|C|}}, \text{ where } \mu_c = \frac{\sum_{i \in I_c} s_i}{|I_c|}.$$

4. Parameter Estimation

In this section, we estimate the parameters in our model. We describe the data in §4.1, and estimate the parameters of the judge schedule module, the defendant arrivals module and the plea bargain module, respectively, in §4.2–4.4.

4.1. Data

We use two datasets from the South Carolina circuit court (i.e., the court of general jurisdiction) between July 1, 2000 and June 30, 2001: the master calendar and the sentencing dataset. The master calendar describes the weekly schedule (i.e., location among the $C = 46$ counties) of the $J = 50$ judges between July 2000 and June 2001. South Carolina trial judges preside over both criminal and civil matters with terms of court typically lasting one week and consisting of only criminal matters (general session) or civil matters (common pleas) for the given term. We restrict our attention in the master calendar to judge-weeks in which the judge is in general session for at least a portion of the week (judges can also, e.g., be on vacation, out sick, in circuit court or involved with the state grand jury), which comprises 854.4 (33.1%) of the 2582 judge-weeks.

The sentencing dataset contains offenders that are convicted of a felony or serious misdemeanor. There are 17,516 sentencing events in the dataset, and 246 (1.4%) of these events were discarded because we could not impute the correct court dates. Of the remaining 17,270 sentencing events (17,012 pleas and 258 trials), 14,977 (14,748 pleas, 229 trials) involve cases where the judge was in general session. Table 2 shows the covariates that we have for each sentencing event, which describe characteristics of the offense and the defendant.

4.2. Judge Schedule Parameters

In this subsection, we compute the weekly judge schedules C_{jct} , the available times in general session T_{jct} and the trial allocations n_{jct}^T , and the mean processing times γ_P and γ_T .

The starting point for our determination of the weekly judge schedules is a mathematical program that assigns judges to counties so as to cover the plea cases from each county while minimizing the degree of judge traveling. In our assignment problem, we focus on satisfying the sentencing cases that were resolved through plea bargains (i.e., excluding those cases that were resolved through trials) in the sentencing dataset. There are two reasons to exclude trials: (i) the primary goal of this work is to understand the impact of the key operational characteristics on plea bargains, which represent 98.5% of the cases in our sentencing dataset, and (ii) trials have a more complex and longer timeline than pleas and the datasets do not give us enough information on the trial schedules.

Table 2 Covariates for each sentencing event

Variable	Description	Values
Offense seriousness	Five-level ordinal score; From the South Carolina crime classification scheme	1 = Misdemeanor 2 = Felony (Class F) 3 = Felony (Class E) 4 = Felony (Class D) 5 = Felony (Class A, B, C or Unclassified)
Commitment score	12-level ordinal measure; Number of commitment offenses	1 = Least serious 12 = Most serious
Offense type	Four-category indicator of the classification of crime committed	1 = Violent 2 = Drug 3 = Property 4 = Other
Mandatory minimum	Minimum prison sentence	1 = Yes 0 = No
Criminal history	Five-level ordinal score; Derived from the number and severity of prior offenses	1 = None 2 = Minimal 3 = Moderate 4 = Considerable 5 = Extensive
Black	Race	1 = Yes 0 = No
Male	Sex	1 = Female 0 = Male
Black \times offense seriousness	Interaction term	$(1, 0) \times (1, 2, 3, 4, 5)$
Black \times crime history	Interaction term	$(1, 0) \times (1, 2, 3, 4, 5)$

Let κ_j be the number of weeks in the year (July 2020 - June 2021) that judge j had a general session assignment in the master calendar (see Table A1 in §EC.1.1 for values), and let d_c equal γ_P (which is estimated below) times the number of pleas processed in general session by county c in the year in the master calendar (see Table A2 in §EC.1.1 for values). Our decision variable is the fraction of judge j 's capacity that is assigned to county c , x_{jc} , and the quadratic function $x_{jc}(1 - x_{jc})$ in (2) forces the optimal x_{jc} values towards 0 or 1, thereby discouraging judge travel. This yields the optimization problem

$$\min \sum_{c=1}^C \sum_{j=1}^J \kappa_j x_{jc} (1 - x_{jc}) \quad (2)$$

$$\text{s.t.} \quad \sum_{j=1}^J \kappa_j x_{jc} \geq d_c, \quad \forall c = 1, \dots, C, \quad (3)$$

$$\sum_{c=1}^C x_{jc} = 1, \quad \forall j = 1, \dots, J, \quad (4)$$

$$0 \leq x_{jc} \leq 1, \quad \forall j = 1, \dots, J, \quad c = 1, \dots, C. \quad (5)$$

The solution x_{jc}^* to this optimization problem is given in Table A3 in §EC.1.1. If $x_{jc}^* > 0$, then we say that county c is one of judge j 's home counties. Because the objective is to minimize the degree of judge traveling, the solution we obtain assigns only six of 50 judges to more than one home county. In Online Supplement §EC.1.1, we provide an algorithm that maps x_{jc}^* into C_{jct} , where $\sum_{c=1}^C C_{jct} = 1$ for all j and t .

To obtain accurate estimates of the mean processing times γ_P and γ_T , we restrict our attention to sentencing events in the general session. Recall T_{jct} be the proportion of week t that judge j is in general session in county c , which can be recovered from the master calendar (this differs from C_{jct} , which is constructed by the algorithm in Online Supplement §EC.1.1) and n_{jct}^T denote the number of trials that judge j handled in county c in week $t = 1, \dots, 52$ in the master calendar. It follows that κ_j in Table A1 in Appendix §A satisfies $\kappa_j = \sum_{c=1}^C \sum_{t=1}^{52} T_{jct}$.

We use linear regression to estimate the mean processing times, γ_P and γ_T . In an analagous manner to above, we define n_{jct}^P to be the number of pleas that judge j handled in county c in week $t = 1, \dots, 52$ in the master calendar. Let $n_{jc}^T = \sum_t n_{jct}^T$ and $n_{jc}^P = \sum_t n_{jct}^P$ be the total number of trials and pleas that judge j handled throughout the year in county c . The regression model, which uses the total number of weeks that judge j worked in county c in the master calendar as the dependent variable, is

$$\sum_{t=1}^{52} T_{jct} = \alpha_c + \gamma_P n_{jc}^P + \gamma_T n_{jc}^T + \epsilon_{jc}, \quad \text{for } j = 1, \dots, J, c = 1, \dots, C,$$

where α_c is interpreted as the average idle time (in weeks) in each county. The regression results (Table EC.2 in the Online Supplement) imply that it takes on average 0.018 weeks to process a plea (i.e., a judge can process $1/0.018=56.2$ pleas per week) and 0.839 weeks to process a trial. With T_{jct} , n_{jct}^T , γ_P and γ_T in hand, we can generate the random variable N_{jct} in §3.2.

As an aside, we note that the total number of judge-weeks in general session during the year is $\sum_{j=1}^J \sum_{c=1}^C \sum_{t=1}^{52} T_{jct} = 854.4$ weeks. In addition, we know from §4.1 that there are 14,748 pleas and 229 trials in general session in the sentencing dataset. Hence, we estimate the judge utilization in the data to be

$$\frac{0.839(229) + 0.018(14,748)}{854.4} = 0.536.$$

4.3. Defendant Arrival Parameters

In this subsection, we compute the weekly arrival rate for each county, λ_c , and determine the covariate values for each arrival, \mathbf{x}_i .

The quantity $\sum_{j=1}^J \kappa_j x_{jc}^*$, where x_{jc}^* is the solution to the mathematical program in (2)-(5), is interpreted as the total number of judge-weeks in a year devoted to county c . Recall that γ_P is the mean processing time (in weeks) for a plea, which is estimated in §4.2. Given a desired value for the utilization ρ , we compute the defendant weekly arrival rate to county c by

$$\lambda_c = \frac{\rho \sum_{j=1}^J x_{jc}^* \kappa_j}{52\gamma_P}, \quad (6)$$

where κ_j appears in Table A2 and x_{jc}^* is given in Table A3 in Appendix §A.

Let D_c denote the set of all defendants from county c in the sentencing dataset and recall that N_c is a Poisson random variable with mean λ_c . Then each week we randomly draw (with replacement) N_c defendants from the set D_c and – following Hester and Hartman (2017) – use the values for the covariates \mathbf{x}_i in Table 2.

4.4. Plea Bargain Parameters

In this subsection, we estimate the defendant probability of incarceration during a plea bargain (p_i), the expected sentence length at trial ($\theta_i \tau_i$), the judge leniency functions ($l_j(\cdot), u_j(\cdot)$), the defendant cost of going to trial ($c_d(i)$) and the defendant delay cost (d).

Defendant Probability of Incarceration During a Plea Bargain. To estimate the probability p_i , we use the covariates \mathbf{x}_i in Table 2 and define the binary variable y_i to equal 1 if defendant i is incarcerated in the sentencing dataset, and to equal 0 otherwise. We fit a logistic regression model using the 17,012 sentencing cases that were resolved by plea bargaining,

$$p_i = \Pr(y_i = 1 | \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})}, \quad (7)$$

where $\boldsymbol{\beta}$ is the vector of regression coefficients for the variables in \mathbf{x}_i . The estimated coefficients are provided in Table EC.3 in the Online Supplement. An alternative regression model incorporated the judge identity as an additional covariate, but the results were very similar and we used the sparser model.

Defendant Expected Sentence Length at Trial. Because θ_i and τ_i appear in our plea bargain module only as the product $\theta_i \tau_i$, it suffices to estimate the product. To predict defendant i 's expected sentence length at trial $\theta_i \tau_i$, we use the hurdle regression model in

Hester and Hartman (2017) and fit the model using the 258 cases that were resolved by trial. The hurdle model uses a logistic regression to predict observations that are zero, and a zero-truncated count model (negative binominal) to predict the remaining nonzero cases:

$$\Pr(\tau_i = 0 | \mathbf{x}_i) = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\omega})}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\omega})} = 1 - \theta_i, \quad (8)$$

$$\Pr(\tau_i | \tau_i > 0, \mathbf{x}_i) = \theta_i \left[\frac{\frac{\Gamma(\tau_i + \alpha^{-1})}{\tau_i! \Gamma(\alpha^{-1})} \left(\frac{\alpha^{-1}}{\alpha^{-1} + \mu} \right)^{\alpha^{-1}} \left(\frac{\mu}{\alpha^{-1} + \mu} \right)^{\tau_i}}{1 - (1 - \alpha\mu)^{-1/\alpha}} \right] \text{ for } \tau_i > 0, \quad (9)$$

where $\boldsymbol{\omega}$ is the set of regression coefficients for the variables in \mathbf{x}_i (Table EC.4 in the Online Supplement), $\Gamma(\cdot)$ is the gamma function, α is a dispersion parameter, and μ is the mean of the negative binomial model (i.e., $\mu = \exp(\mathbf{x}_i^T \boldsymbol{\gamma})$, where $\boldsymbol{\gamma}$ is the set of regression coefficients for the negative binomial regression model; see Table EC.5 in the Online Supplement).

For defendant i , we predict his probability of conviction at trial, $\hat{\theta}_i$, using (8) and the coefficients in Table EC.4 in the Online Supplement, and predict his sentence length upon conviction at trial, $\hat{\tau}_i$, using (9) and the coefficients in Table EC.5 in the Online Supplement. The expected trial sentence length of defendant i is estimated to be $\hat{\theta}_i \hat{\tau}_i$.

Judge Leniency Functions. To estimate the judge leniency thresholds $\ell_j(\cdot)$ and $u_j(\cdot)$ for each judge as functions of the expected sentence length at trial, $\theta_i \tau_i$, we restrict attention in the sentencing dataset to those defendants who accepted a plea bargain and were incarcerated. In the sentencing dataset, let s_i be defendant i 's sentence and \mathcal{I}_j be the set of plea bargains handled by judge j for $j = 1, \dots, J$.

To estimate the leniency functions for a particular judge j , we create a scatter plot of $(\hat{\theta}_i \hat{\tau}_i, s_i)$ for $i \in \mathcal{I}_j$. The key idea is to construct a convex hull for judge j generated by the origin $(0, 0)$ and the points $(\hat{\theta}_i \hat{\tau}_i, s_i)$ for $i \in \mathcal{I}_j$, and use the upper and lower boundary of the convex hull (given by the two red dots in Figure 2) to estimate the maximum and the minimum sentence length that judge j would approve for defendant i with a mean sentence length at trial equal to $\hat{\theta}_i \hat{\tau}_i$ (i.e., values of $u_j(\theta_i \tau_i)$ and $\ell_j(\theta_i \tau_i)$).

We modify this approach in two ways that are described in Online Supplement §EC.1.2. First we detect and remove outliers from the observations $\{(\hat{\theta}_i \hat{\tau}_i, s_i) | s_i > 0, i \in \mathcal{I}_j\}$ using the Mahalanobis Distance, which results in 3.44% of observations being identified as outliers, with a range from 0% to 5.77% for each judge. Second, so that we can estimate the judge leniency thresholds for defendants whose expected sentence length at trial is greater

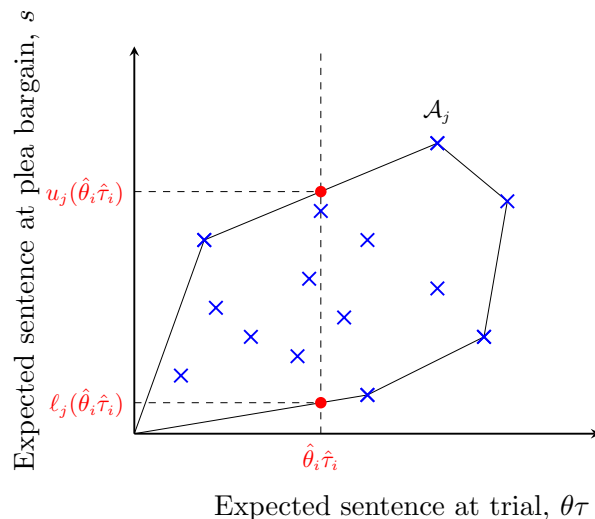


Figure 2 A scatter plot of the sentencing handled by judge j and the constructed convex hull \mathcal{A}_j . This figure is for illustration purposes and is not based on the data.

than $\max_{i \in I_j} \hat{\theta}_i \hat{\tau}_i$, we extrapolate the convex hull by including two artificial points that are determined based on the maximum possible expected trial sentence obtained from the sentencing dataset. After these two modifications, the resulting convex hulls for all judges appear in Fig. EC.4 in the Online Supplement.

To provide a sense of how these leniency functions impact the sentences in our model, we use each judge's convex hull to compute the mean sentence for all defendants in the sentencing dataset that received a nonzero sentence (Fig. EC.5 in the Online Supplement). The mean overall sentence length is 51.6 months, and the mean for individual judges ranges from 44.6 (Judge 36) to 62.0 (Judge 15).

Defendant Trial Cost. To estimate the cost of going to trial, $c_d(i)$, we let $j(i)$ denote the presiding judge for defendant i 's case in the sentencing dataset. We have four cases to consider for defendant i : (1) he received sentence $s_i > 0$ in a plea bargain that satisfies $\ell_{j(i)}(\hat{\theta}_i \hat{\tau}_i) < s_i < u_{j(i)}(\hat{\theta}_i \hat{\tau}_i)$, (2) he received a sentence $s_i = u_{j(i)}(\hat{\theta}_i \hat{\tau}_i)$ in a plea bargain, (3) he received his sentence through a plea bargain and was not incarcerated, i.e., $s_i = 0$, and (4) his case went to trial.

We partition all cases according to these four outcomes by defining the sets

$$\begin{aligned} \mathcal{I}^1 &= \cup_{j=1}^J \{i \in \mathcal{I}_j \mid s_i > 0, \text{ and } \ell_{j(i)}(\hat{\theta}_i \hat{\tau}_i) < s_i < u_{j(i)}(\hat{\theta}_i \hat{\tau}_i)\}, \\ \mathcal{I}^2 &= \cup_{j=1}^J \{i \in \mathcal{I}_j \mid s_i > 0, \text{ and } s_i = u_{j(i)}(\hat{\theta}_i \hat{\tau}_i)\}, \end{aligned}$$

$$\mathcal{I}^3 = \cup_{j=1}^J \{i \in \mathcal{I}_j | s_i = 0\},$$

$$\mathcal{I}^4 = \cup_{j=1}^J (\mathcal{I}_j \setminus \{\mathcal{I}_j^1 \cup \mathcal{I}_j^2 \cup \mathcal{I}_j^3\}).$$

We estimate the trial cost $c_d(i)$ when i belongs to $\mathcal{I}^1, \mathcal{I}^2, \mathcal{I}^3$, and \mathcal{I}^4 separately. When $i \in \mathcal{I}^1$, our plea bargain model implies that

$$s_i = \min \left\{ \hat{\theta}_i \hat{\tau}_i + c_d(i), u_{j(i)}(\hat{\theta}_i \hat{\tau}_i) \right\}, \quad (10)$$

$$= \hat{\theta}_i \hat{\tau}_i + c_d(i) \quad \text{because } s_i < u_{j(i)}(\hat{\theta}_i \hat{\tau}_i). \quad (11)$$

Thus, we can infer defendant i 's cost of trial as $c_d(i) = s_i - \hat{\theta}_i \hat{\tau}_i$ in this case.

For the other three cases $l = 2, 3, 4$, we define a set S^l of similar defendants, and use a k -nearest neighbor algorithm with $k = 20$ (see Online Supplement §EC.1.3 for a discussion of how we set k) to estimate the trial cost for defendants in groups $\mathcal{I}^2, \mathcal{I}^3$ and \mathcal{I}^4 . More specifically, for $l = 2, 3, 4$, if $|S^l| < 20$ the estimated trial cost is the average value of $c_d(i')$ where $i' \in S^l$. If $|S^l| \geq 20$, then we select 20 observations that are the nearest neighbors to observation $(\hat{\theta}_i \hat{\tau}_i, s_i)$ using the Mahalanobis distance from $\{(\hat{\theta}_i \hat{\tau}_i, s_i)\} \cup S^l$. We let S^l denote the set of those selected observations and estimate the trial cost by the average value of $c_d(i')$ where $i' \in S^l$.

It remains to compute the sets S^2, S^3 and S^4 . When $i \in \mathcal{I}^2$, we can only infer that $c_d(i) \geq u_{j(i)}(\hat{\theta}_i \hat{\tau}_i) - \hat{\theta}_i \hat{\tau}_i$. Because set \mathcal{I}^1 contains most of the defendants with a positive sentence and we can infer their trial costs, we use those estimates to infer the trial costs of the defendants in set \mathcal{I}^2 . Hence we let

$$S^2 = \left\{ i' \in \mathcal{I}^1 : c_d(i') \geq u_{j(i)}(\hat{\theta}_i \hat{\tau}_i) - \hat{\theta}_i \hat{\tau}_i \right\}.$$

When $i \in \mathcal{I}^3$, we cannot infer any bounds on $c_d(i)$, and therefore let $S^3 = \mathcal{I}^1$. When $i \in \mathcal{I}^4$, we infer that $c_d(i) < \ell_{j(i)}(\hat{\theta}_i \hat{\tau}_i) - \hat{\theta}_i \hat{\tau}_i$, and let

$$S^4 = \left\{ i' \in \mathcal{I}^1 : c_d(i') \leq \ell_{j(i)}(\hat{\theta}_i \hat{\tau}_i) - \hat{\theta}_i \hat{\tau}_i \right\}.$$

The histograms of $c_d(i)$ for $i \in \mathcal{I}^i, i = 1, \dots, 4$ appear in Fig. EC.6 in the Online Supplement, where the cost of going to trial is estimated to be negative for many defendants. We interpret this to mean that many defendants should prefer to go to trial rather than to plea bargain.

Defendant Waiting Cost. To estimate d , we simulate ten 50-year replications of the system, and collect performance measures for the middle 30 years of each replication. Because we use the actual judge schedule in the master calendar, the operational characteristics η and ρ do not apply. We set the shopping window to $r = 4$ weeks, which is the base-case value in §5. We use the 14,977 sentencing events in general session to estimate the weekly arrival rate and allow it to vary across months to more accurately mimic the seasonality in the data, and randomize the timing and characteristics of the defendant arrivals as described earlier. Because τ_i and $c_d(i)$ are measured in months whereas $w_i(j)$ is measured in weeks in (1), the cost d is in units of months per week. Consequently, with an assumption of four weeks per month, $d = 0.25$ means that the cost of delaying a court case for one week is equal to one week of detention. We consider $\{0, 0.1, 0.25\}$ as possible values for d . The three performance measures are quite insensitive to d (Table 3) and we use $d = 0.1$.

Table 3 Statistics from the sentencing dataset and simulation results with weekly waiting cost $d \in \{0, 0.1, 0.25\}$ and a defendant shopping window $r = 4$

	Data	$d = 0$	$d = 0.1$	$d = 0.25$
Mean plea sentence (for nonzero sentences)	54.27	50.21	50.29	50.29
Standard deviation of plea sentence	55.18	52.36	52.34	52.34
Percentage of trial cases	1.51%	1.94%	1.99%	1.99%
Standard deviation across counties (for nonzero sentences)	11.74	12.52	12.50	12.50

5. Results

We investigate the impact of the key operational characteristics in §5.1, and in §5.2 consider a hypothetical urban model in which all defendants and judges are located in one place (i.e., all defendants have access to all judges). All simulation results consider 10 independent replications of 50 years, with the first and last 10 years of each replication discarded.

5.1. Impact of the Key Operational Characteristics

In this subsection, we investigate the impact of the judge travel probability η , the defendant shopping window r and the judge utilization ρ on the three performance measures. We consider base values of $\eta = 0.5$, $r = 4$ weeks (although South Carolina published annual schedules for judges, these schedules were updated frequently) and $\rho = 0.5$. Figs. 3-8 show the performance measures (the shaded areas surrounding the

curves are 95% confidence intervals) plotted against one of the operational characteristics (with values $\eta \in \{0, 0.1, 0.2, \dots, 0.9\}$, $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks and $\rho \in \{0.5, 0.7, 0.9, 0.95, 0.99\}$), with a second operational characteristic set at its base value, and several different curves for changes in the third operational characteristic. In addition, three-dimensional plots of the performances measures versus two operational characteristics with the third set at its base value appear in Figs. EC.7-EC.9 in the Online Supplement.

We defer a discussion of the county variation (Figs. 3c-8c) and focus on the impact on the mean and standard deviation of the plea sentences. As the travel probability η increases (Figs. 3-4), the mean and standard deviation of the plea sentences decreases with diminishing impact (i.e., the curves are convex), and the impact is larger when the shopping window r is large or the judge utilization ρ is small. Judges travel more as η increases, increasing defendants' shopping opportunities and therefore their likelihood of securing a more lenient judge. A larger shopping window gives defendants a longer period to choose their preferred judges. However, when judge utilization is high, defendants cannot always choose a more lenient judge because the lenient judges may be overwhelmed with other defendants' pleas.

Similarly, as the shopping window r increases (Figs. 5-6), the mean and standard deviation of the plea sentence decrease with decreasing returns, with the impact being larger when the travel probability η is large or the judge utilization ρ is small.

In contrast, when the judge utilization ρ increases (Figs. 7-8), the mean and standard deviation of the plea sentence increase, and are slightly convex (i.e., exhibit increasing returns), with the impact being greater for larger values of the travel probability η or the shopping window r . Overall, at least over the ranges considered here, the impact of judge utilization is smaller than the impacts of the travel probability and the shopping window.

Returning to the impact of the key operational characteristics on the standard deviation of the mean plea sentence across counties, we see that the county variation typically (with a few exceptions) decreases as the travel probability η or the judge utilization ρ increases (Figs. 3c-4c, 7c-8c); in the latter case, a high utilization reduces the chance of choosing a lenient judge for defendants across all counties rather than benefiting some counties more.

However, the behavior of county variation is more complex when the shopping window r increases, both with a fixed ρ or a fixed η . When we fix the judge utilization at $\rho = 0.5$ (Fig. 5c), if the travel probability η is small then increasing the shopping window r

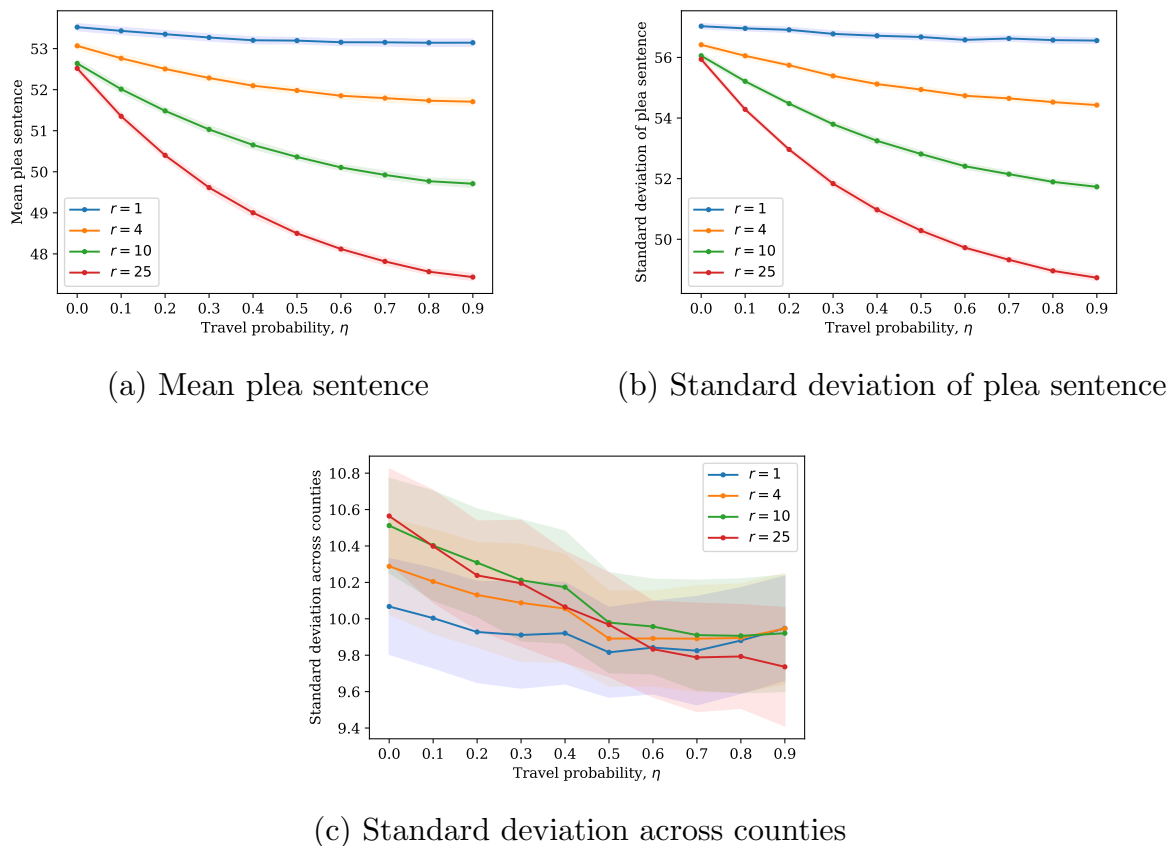
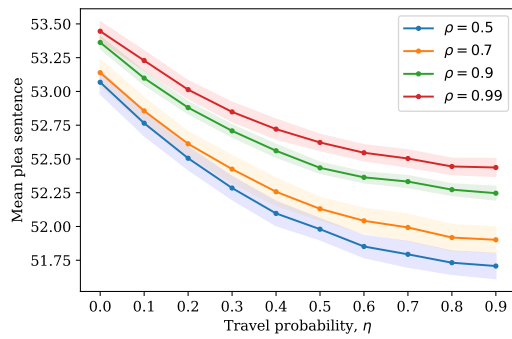


Figure 3 Performance measures versus the travel probability $\eta \in \{0, 0.1, \dots, 0.9\}$, where the defendant shopping window is $r \in \{1, 4, 10, 25\}$ weeks and judge utilization is fixed at $\rho = 0.5$.

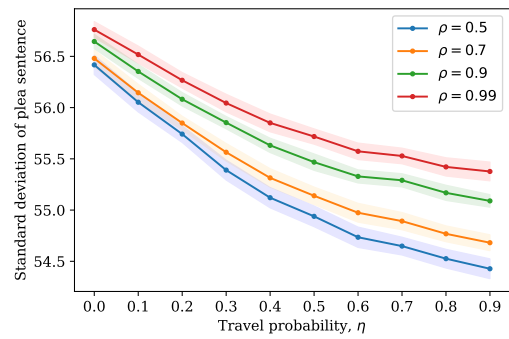
primarily impacts the counties with multiple home judges, which initially increases the county variation. In contrast, when η is large, an increase in the shopping window allows defendants in all counties to shop, leading to a decrease in county variation.

When we fix the travel probability at $\eta = 0.5$ (Fig. 6c), an increasing shopping window r enlarges the defendants' shopping options. A longer shopping window provides a larger benefit to defendants in counties with more visiting judges. When r is large enough, however, the defendants from most counties are able to see all possible judges from their shopping window, which stabilizes the county variation.

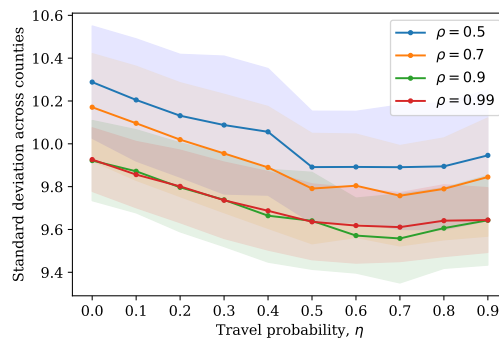
We perform three more calculations in this subsection in an attempt to gain more insight. First, to put the range of performance in Figs. 3-8 in perspective, we simulate an idealized scenario where defendants have access to all judges ($\eta = 1, r = \infty$), and judges have infinite capacity to handle pleas ($\rho = 0$). The results of this idealized scenario (Table 4) suggest



(a) Mean plea sentence



(b) Standard deviation of plea sentence



(c) Standard deviation across counties

Figure 4 Performance measures versus the travel probability $\eta \in \{0, 0.1, \dots, 0.9\}$, where the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ and the shopping window is fixed at $r = 4$ weeks.

that large values of the travel probability and the shopping window (e.g., Figs. 3-5) can achieve over half of the impact of the idealized scenario.

Table 4 Results from an idealized scenario ($\eta = 1$, $r = \infty$, $\rho = 0$)

Mean plea sentence	43.75
Standard deviation of plea sentence	49.28
County variation	9.12

Second, referring back to the three possible outcomes in the plea bargain model in §3.3, we simulate the model under the base-case values ($\eta = 0.5$, $r = 4$, $\rho = 0.5$) and find that 6.4% of cases go to trial, the plea sentence equals the judge’s upper allowable sentence in 9.2% of the cases, and the plea sentence is equal to the defendant’s expected total trial cost in 84.5% of the cases.

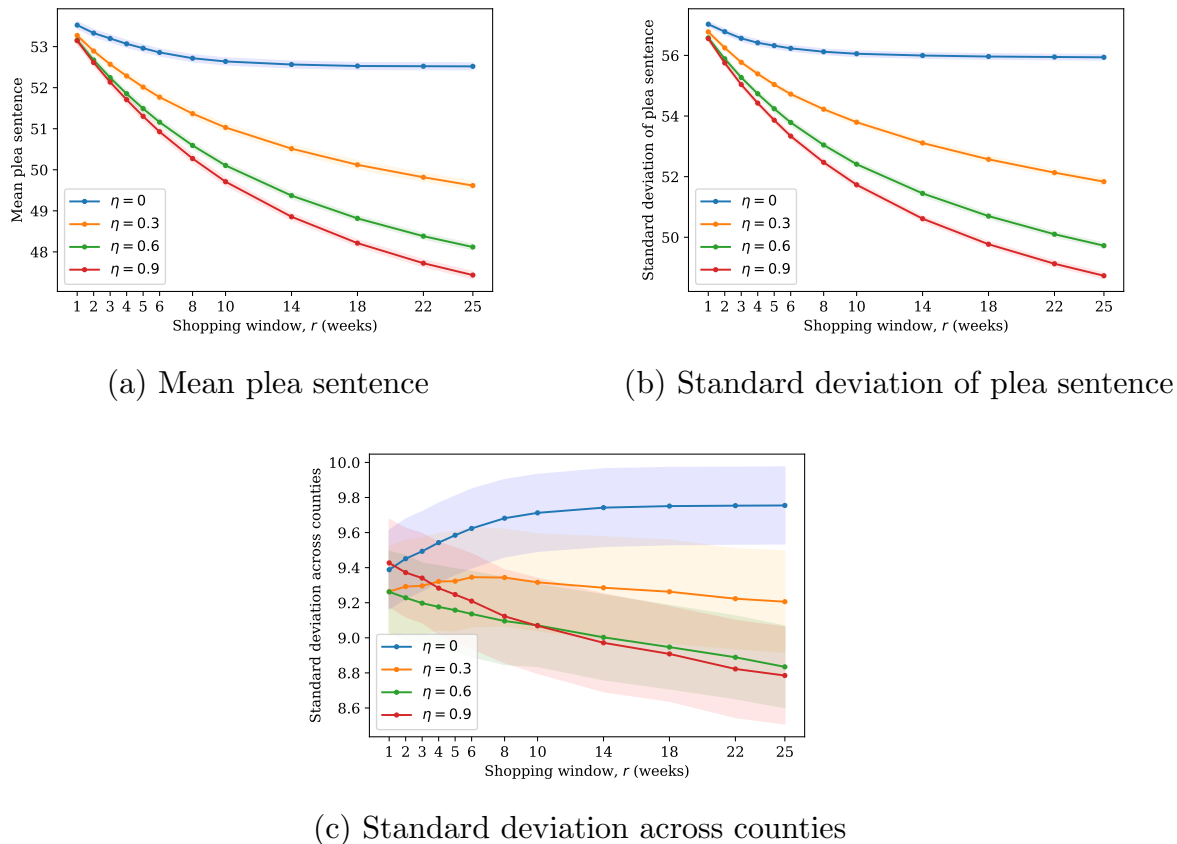


Figure 5 Performance measures versus the shopping window $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks, where the travel probability $\eta \in \{0, 0.3, 0.6, 0.9\}$ and judge utilization is fixed at $\rho = 0.5$.

Finally, we elaborate on (and provide a visualization of) this last calculation by roughly estimating an upper bound on the proportion of defendants who are impacted by judge shopping. We consider a hypothetical two-judge scenario in which each defendant is assigned the most lenient judge (Judge 36 in Fig. EC.5 in the Online Supplement) with probability 0.60 and the harshest judge (Judge 15 in Fig. EC.5 in the Appendix) with probability 0.40, where these probabilities maintain the mean nonzero sentence length to be 51.6 months, which is the overall mean sentence length in Fig. EC.5 in the Online Supplement. The defendants whose $(\theta_i \tau_i, s_i)$ fall in the interior of both convex hulls in Fig. EC.10 in the Online Supplement, which accounts for 90.5% of the defendants, are indifferent between Judges 15 and 36 because they would receive the same sentence length, $\theta_i \tau_i + c_d(i)$, from either judge. In contrast, the defendants with higher $\theta_i \tau_i + c_d(i)$ are more likely to be sensitive to the choice of judge. These last two calculations suggest that a small

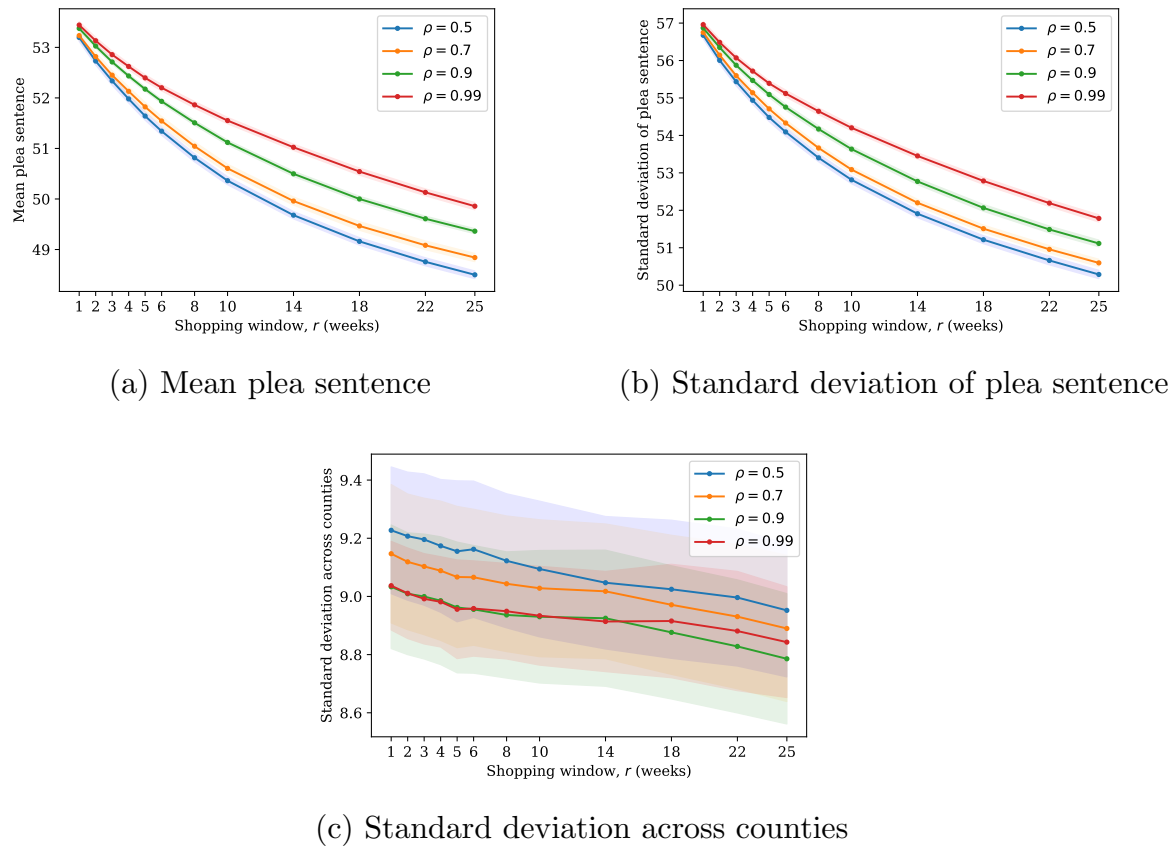
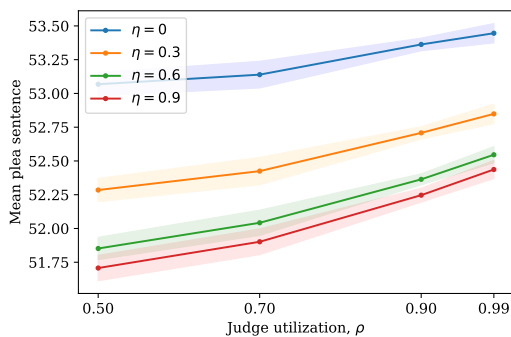


Figure 6 Performance measures versus the shopping window $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks, where the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ and the travel probability is fixed at $\eta = 0.5$.

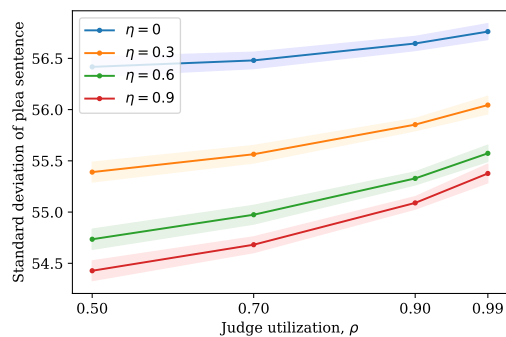
proportion of defendants are impacted, but that the impact on some of these individuals is large.

5.2. Urban Model

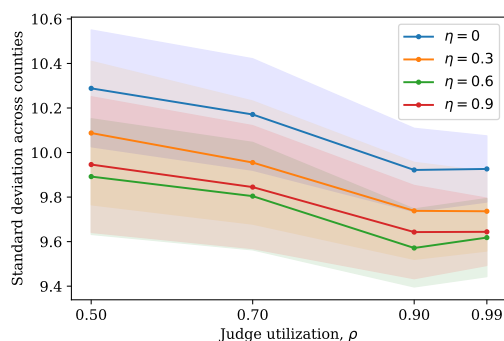
In this subsection, we adapt our model from the South Carolina setting to an urban setting, where all defendants and judges are located in the same place. We assume that defendants have access to all judges. This assumption is not realistic in that a case would generally be assigned to a judge without input from a defendant. However, local rules and practices may allow parties to engage in some strategic maneuvering, and our analysis is intended to gain insight into the impact of such maneuvering. The traveling probability η now plays no role, and we introduce another parameter whose value is specified: k , which represents the number of judges. Our simulation model is the same as before, except we now allow all $J = 50$ judges to be accessible to each defendant. Recall that an arriving defendant in our simulation model is assigned a set $J_i(r)$ of judges assigned to the defendant's county with



(a) Mean plea sentence



(b) Standard deviation of plea sentence



(c) Standard deviation across counties

Figure 7 Performance measures versus the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.95, 0.99\}$, where the travel probability $\eta \in \{0, 0.3, 0.6, 0.9\}$ and the shopping window is fixed at $r = 4$ weeks.

remaining capacity within r weeks, and then chooses a judge according to (1). We change the model by allowing $J_i(r)$ to include all $J = 50$ judges that have remaining capacity within r weeks of the defendant's arrival. If the number of judges in this set, $|J_i(r)|$, satisfies $|J_i(r)| \leq k$, then the defendant chooses the best among the $|J_i(r)|$ judges according to (1). If $|J_i(r)| > k$, then the simulation chooses k out of $|J_i(r)|$ judges at random, and the defendant chooses the best among these k judges according to (1). Hence, the system behaves as if there are k judges, but this construction allows us to use the characteristics of all $J = 50$ judges in the dataset.

As in §5.1, we plot our performances measures (county variation no longer plays a role) against two key operational characteristics, keeping the third one fixed (Figs. 9-14, with the three-dimensional plots displayed in Figs. EC.11-EC.13 in the Online Supplement), considering $k \in \{1, 2, 3, 4, 5, 6, 7, 10, 15, 20\}$ and a base-case value of $k = 4$ judges.

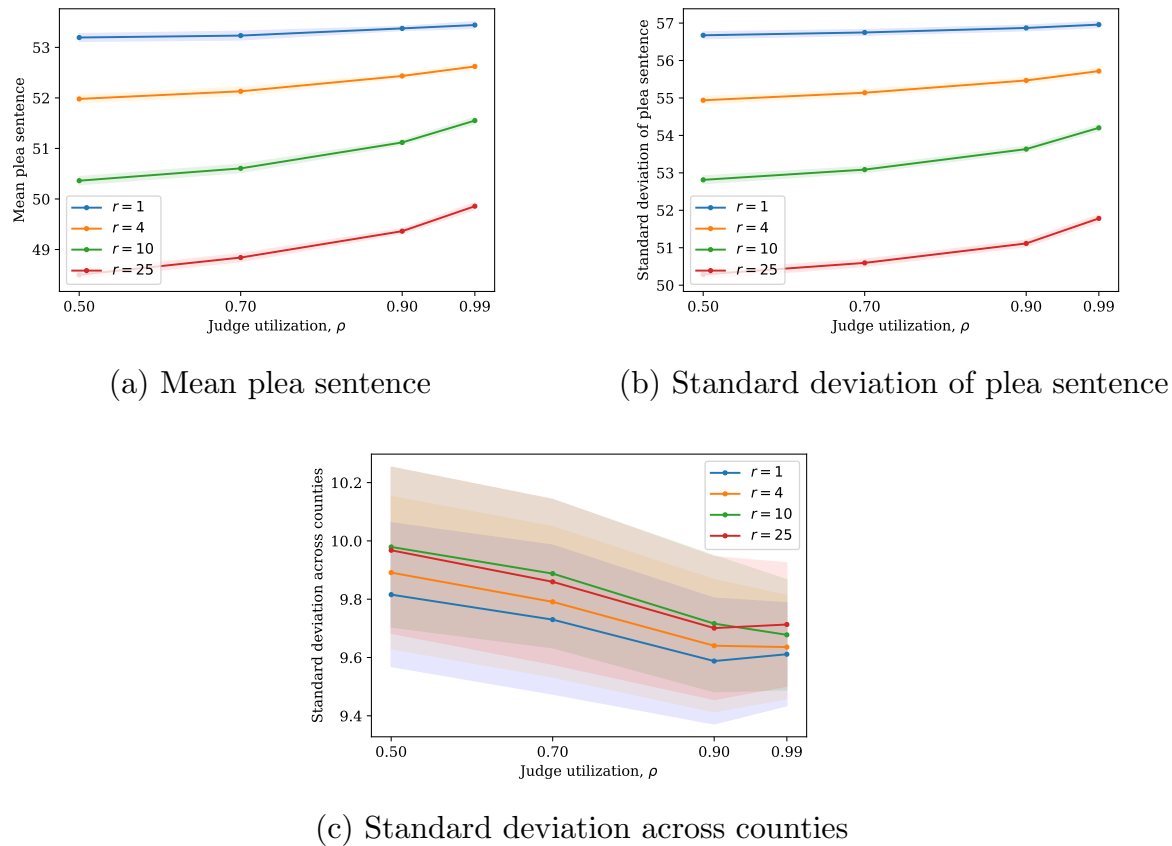
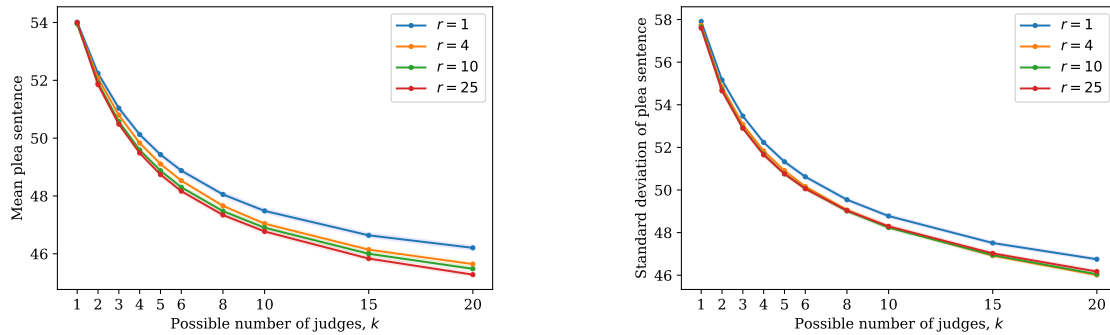


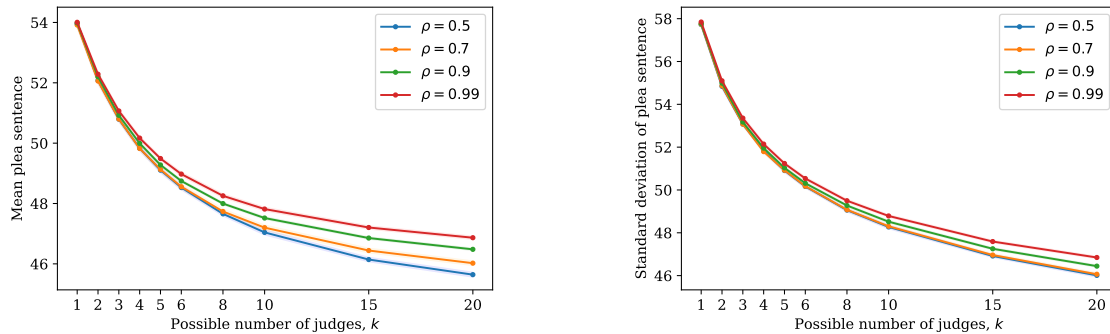
Figure 8 Performance measures versus the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.95, 0.99\}$, where the shopping window $r \in \{1, 4, 10, 25\}$ weeks and the travel probability is fixed at $\eta = 0.5$.

We highlight the differences in the results between the urban model and the South Carolina model. As expected, the mean and standard deviation of the sentence length decreases as the number of judges k increases because increasing k provides more shopping options. While this impact increases with a larger shopping window r and a smaller judge utilization ρ , these dependencies (Figs. 9-10) are much smaller than in the corresponding results (Fig. 3) in the South Carolina model. Most notably, a shopping window of $r = 1$ week still allows defendants shopping opportunities when k is large. In addition, a comparison of Figs. 6 and 11 reveals that the system performance is very insensitive to the judge utilization ρ in the urban model when the shopping window r is large, because the defendant gets more choices in the urban model than the South Carolina model when ρ and r are large. In Fig. 13, the mean plea sentence is insensitive to the judge utilization ρ when the shopping window r is large because most defendants have access to k judges when r is large, regardless of the utilization.



(a) Mean plea sentence (b) Standard deviation of plea sentence

Figure 9 Performance measures versus the number of judges $k \in \{1, 2, 3, 4, 5, 6, 8, 10, 15, 20\}$ in the urban model, where the shopping window $r \in \{1, 4, 10, 25\}$ weeks and the judge utilization is fixed at $\rho = 0.5$.

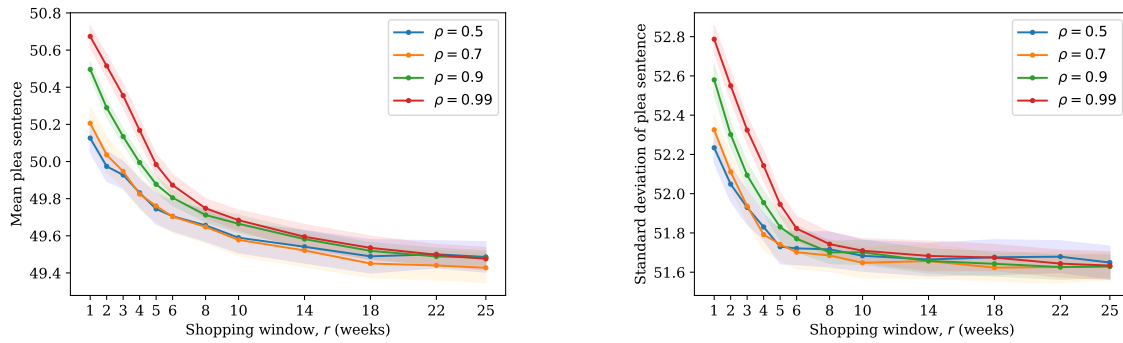


(a) Mean plea sentence (b) Standard deviation of plea sentence

Figure 10 Performance measures versus the number of judges $k \in \{1, 2, 3, 4, 5, 6, 8, 10, 15, 20\}$ in the urban model, where the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ and the shopping window is fixed at $r = 4$ weeks.

6. Discussion

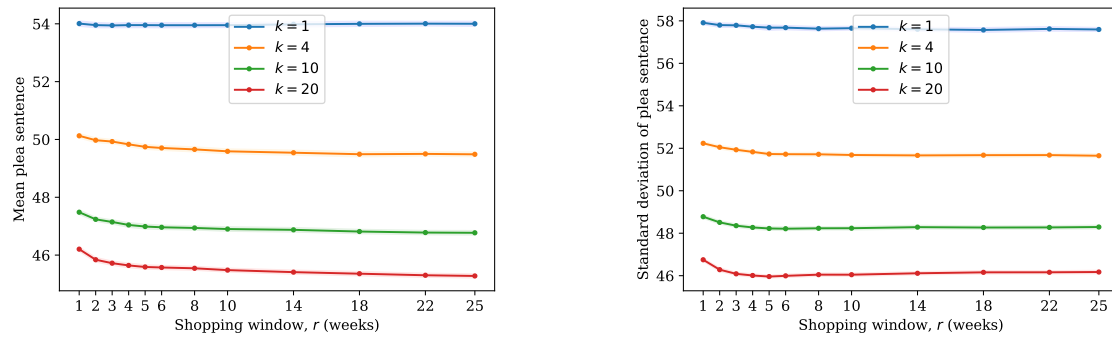
Motivated by the insights and data in Hester (2017), we formulate and calibrate a mathematical model that allows judges to rotate among counties, allows defendants to shop for lenient judges, and incorporates a sentencing model that captures the strategic interactions among a judge, a defendant and a prosecutor. Our goal is to gain an understanding of how three key operational characteristics – the amount of judicial rotation (as measured by the judge travel probability), the amount of leeway defendants have in shopping (as measured by the shopping time window) and the system congestion (as measured by the judge utilization) – impact the mean and standard deviation of nonzero sentence lengths, and the county variation in mean sentence lengths. The latter two of these three performance measures quantify the amount of sentencing inequity across defendants.



(a) Mean plea sentence

(b) Standard deviation of plea sentence

Figure 11 Performance measures versus the shopping window $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks in the urban model, where the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ and the number of judges is fixed at $k = 4$.

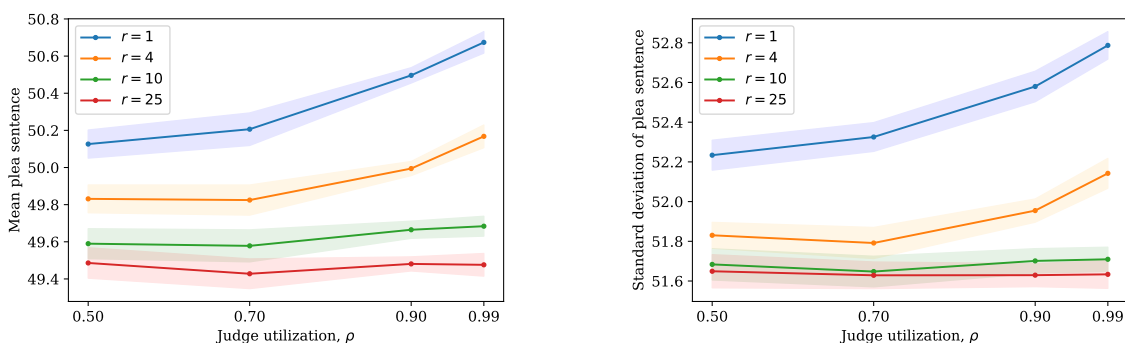


(a) Mean plea sentence

(b) Standard deviation of plea sentence

Figure 12 Performance measures versus the shopping window $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks in the urban model, where the number of judges $k \in \{1, 2, 3, 4, 5, 6, 8, 10, 15, 20\}$ and the judge utilization is fixed at $\rho = 0.5$.

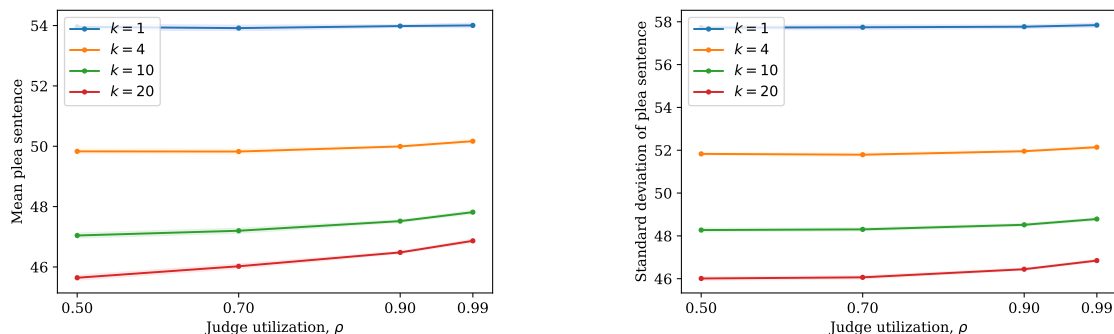
Our first-order findings (Figs. 3-8) are intuitive: the mean and standard deviation of the plea sentence length decrease when judges travel more among counties, defendants have a longer shopping window, and judges have more excess capacity. All three characteristics make it easier for a defendant to plea in front of a lenient judge. The impact of all three characteristics have decreasing returns to scale, in that a moderate amount of judge shopping, a moderate shopping window, and a moderate amount of excess judicial capacity achieve a significant proportion of the potential effects. In terms of interaction effects, we find that the judge travel probability and the shopping window are complements: each variable generates a bigger impact when the other variable is large. That is, very little



(a) Mean plea sentence

(b) Standard deviation of plea sentence

Figure 13 Performance measures versus the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ in the urban model, where the shopping window $r \in \{1, 2, 3, 4, 5, 6, 8, 10, 14, 18, 22, 25\}$ weeks and the number of judges is fixed at $k = 4$.



(a) Mean plea sentence

(b) Standard deviation of plea sentence

Figure 14 Performance measures versus the judge utilization $\rho \in \{0.5, 0.7, 0.9, 0.99\}$ in the urban model, where the number of judges $k \in \{1, 4, 10, 20\}$ and the shopping window is fixed at $r = 4$ weeks.

impact occurs when judges travel but defendants cannot shop, or when defendants can shop but the judges do not travel. In contrast, high judge utilization reduces – but does not negate – the first-order and synergistic effects of increasing the judge travel probability and the shopping window. In addition, although the overall reduction on the mean and standard deviation is modest ($\approx 10\%$), we find that the changes in sentencing as a result of these key operational characteristics affect a small proportion of defendants – typically those who have a high cost of going to trial – but the impact on these individual defendants is large.

The Hester (2017) study uses data from South Carolina in 2000-2001, when judicial rotation was used and most counties had a single home judge Hester (2017). Very few

states have employed judicial rotation during the last 20 years. In contrast, urban areas throughout the US typically have multiple judges in the same county, which may allow for some judge shopping by defendants. Hence, we adapt the South Carolina model into an urban model, where the number of judges replaces the judge travel probability as the third key operational characteristic. As expected, the mean and standard deviation of the plea sentence length are reduced when the number of judges and the shopping window are large and the judge utilization is low. However, in contrast to the county model, the impact of the number of judges is much larger than the impact of the shopping window, and there is very little synergy: even if intertemporal shopping is not allowed (i.e., $r = 1$ week), the defendant still gets to choose the most lenient among the available judges and can lower his mean sentence length. The impact of judge utilization is smaller than it is in the South Carolina model.

The system we attempt to model is very complex, and our model should be viewed as a mere caricature of the actual process. In particular, we use a highly idealized model for the plea bargain process (e.g., we do not include costs of going to trial for the prosecutor, as in Silveira (2017)), and a simplified model for how judges spend their time and process cases (e.g., judicial decisions do not depend on the amount of congestion, as in Yang (2016)). Moreover, there are no data to directly estimate some of the parameters, particularly the cost of going to trial and the waiting cost. However, some of our methods may be useful, particularly our construction of the convex hulls to estimate judge leniency.

These limitations lead us to conclude that the exact numerical results in Figs. 3-8 should not be interpreted as accurate counterfactual predictions of the system behavior in South Carolina in 2000-2001. However, we believe that our broad qualitative insights are likely to be robust. For one, in relation to previous work in South Carolina, a jurisdiction with heavy utilization of judicial rotation, our findings articulate distinct dimensions of rotation operationalized as the amount of rotation, the shopping window, and capacity utilization. The findings suggest that increased rotation and a longer shopping window can lead to larger decreases in the mean and standard deviation of sentence length, with a small impact of capacity utilization. Thus, expanding choice of judges could assist in aims of incarceration reduction and achieving greater uniformity in sentencing. Second, as applied in the urban models, introducing greater choice of judge could advance similar aims, even with short shopping windows. We do note the paradox, however, that introducing greater

shopping choice in an urban jurisdiction may facilitate inter-judge uniformity within the urban jurisdiction but thereby exacerbate county differences within a state. As a result, some defendants in rural settings may suffer in comparison if their county has a single harsh judge.

References

- Eisenstein J, Flemming RB, Nardulli PF (1988) *The Contours of Justice: Communities and Their Courts*. (Boston, MA: Little, Brown).
- Frazier CE, Bock EW (1982) Effects of court officials on sentence severity: Do judges make a difference? *Criminology* 20(2):257–272.
- Gross SR, Syverud KD (1991) Getting to no: A study of settlement negotiations and the selection of cases for trial. *Michigan Law Review* 90:319–393.
- Hester R (2017) Judicial rotation as centripetal force: Sentencing in the court communities of South Carolina. *Criminology* 55(1):205–235.
- Hester R, Hartman TK (2017) Conditional race disparities in criminal sentencing: A test of the liberation hypothesis from a non-guidelines state. *Journal of Quantitative Criminology* 33(1):77–100.
- Johnson BD (2006) The multilevel context of criminal sentencing: Integrating judge- and county-level influences*. *Criminology* 44(2):259–298.
- Myers MA, Talarico SM (1987) *The Social Contexts of Criminal Sentencing* (Springer New York, NY).
- Ostrom B, Hanson R, for State Courts NC, Institute APR (1999) *Efficiency, Timeliness, and Quality: A New Perspective from Nine State Criminal Trial Courts* (National Center for State Courts).
- Silveira BS (2017) Bargaining with asymmetric information: An empirical study of plea negotiations. *Econometrica* 85(2):419–452.
- Spohn C (1990) The sentencing decisions of black and white judges: Expected and unexpected similarities. *Law & Society Review* 24(5):1197–1216.
- Steffensmeier D, Britt CL (2001) Judges' race and judicial decision making: Do black judges sentence differently? *Social Science Quarterly* 82(4):749–764.
- Ulmer JT (1997) *Social Worlds of Sentencing: Court Communities Under Sentencing Guidelines* (SUNY Press, NY).
- Ulmer JT (2019) Criminal courts as inhabited institutions: Making sense of difference and similarity in sentencing. *Crime and Justice* 48:483–522.
- Wang C (2019) *Three Data-driven Model-based Policy Analyses in Criminology*. Ph.D. thesis, Stanford University.
- Yang CS (2016) Resource constraints and the criminal justice system: Evidence from judicial vacancies. *American Economic Journal: Economic Policy* 8(4):289–332.

Appendix A: Additional Results for Judge Schedule Parameter Estimation

In this section, we provide the required values κ_j , where $j \in J$, and d_c , where $c \in C$, that we use for solving the optimization problem shown in Equations (2)-(5) in Tables A1 and A2, respectively. We also provide the solution to the optimization problem in Table A3.

Table A1 Values for κ_j

Judge Number	κ_j	Judge Number	κ_j	Judge Number	κ_j	Judge Number	κ_j
1	25.0	14	15.0	28	19.6	41	8.0
2	24.4	15	11.0	29	20.4	42	9.0
3	15.0	17	14.6	30	20.2	43	14.6
4	17.5	18	25.4	31	12.0	44	16.0
5	15.4	19	19.0	32	13.6	45	18.2
6	18.2	20	4.8	33	20.0	46	10.33
7	34.0	21	14.2	34	15.2	47	20.2
8	14.8	22	20.7	35	11.0	48	16.0
9	21.3	23	13.2	36	10.0	49	27.8
10	22.2	24	22.7	37	9.5	50	26.3
11	22.0	25	17.8	38	13.0		
12	16.8	26	21.2	39	15.8		
13	21.9	27	12.6	40	8.0		

Table A2 Values for d_c

County Number	d_c	County Number	d_c	County Number	d_c	County Number	d_c
1	13.581	13	6.835	25	2.617	37	0.837
2	2.812	14	3.969	26	10.644	38	2.083
3	2.100	15	4.058	27	11.463	39	2.848
4	4.788	16	23.211	28	0.481	40	2.599
5	1.958	17	2.011	29	6.408	41	1.050
6	5.910	18	1.780	30	28.569	42	1.299
7	1.691	19	4.343	31	0.854	43	5.589
8	2.617	20	27.270	32	1.2816	44	1.371
9	4.948	21	3.097	33	11.606	45	0.374
10	0.552	22	2.812	34	13.635	46	5.020
11	2.207	23	19.010	35	1.086		
12	6.177	24	2.741	36	3.115		

Table A3 Judge assignment obtained by solving (2)-(5)

Judge Number	Assigned counties (Percentage)
1	County 34 (100%)
2	County 37 (100%)
3	County 46 (100%)
4	County 13 (100%)
5	County 35 (100%)
6	County 45 (4.0%), County 12 (58.0%), County 1 (38.0%)
7	County 2 (100%)
8	County 6 (100%)
9	County 20 (100%)
10	County 44 (100%)
11	County 38 (100%)
12	County 21 (100%)
13	County 31 (100%)
14	County 23 (100%)
15	County 5 (100%)
16	County 35 (100%)
17	County 14 (100%)
18	County 16 (100%)
19	County 13 (100%)
20	County 4 (100%)
21	County 36 (100%)
22	County 43 (100%)
23	County 33 (100%)
24	County 36 (100%)
25	County 38 (100%)
26	County 16 (100%)
27	County 10 (100%)
28	County 30 (100%)
29	County 33 (100%)
30	County 39 (17.0%), County 26 (5.0%), County 24 (29.0%), County 25 (25.0%), County 40 (25.0%)
31	County 35 (100%)
32	County 3 (23.0%), County 20 (14.0%), County 41 (12.0%), County 4 (1.0%), County 9 (50.0%)
33	County 27 (100%)
34	County 30 (100%)
35	County 13 (100%)
36	County 20 (100%)
37	County 27 (100%)
38	County 1 (100%)
39	County 26 (100%)
40	County 39 (40.0%), County 17 (60.0%)
41	County 15 (100%)
42	County 42 (21.0%), County 4 (0.0%), County 26 (1.0%), County 8 (55.0%), County 32 (22.0%)
43	County 11 (100%)
44	County 29 (53.0%), County 26 (13.0%), County 19 (34.0%)
45	County 18 (100%)
46	County 22 (100%)
47	County 7 (100%)
48	County 23 (100%)
49	County 28 (100%)
50	County 11 (100%)