

---

# Outliers Exist: What Happens if You are a Data-Driven Exception?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Data-driven tools are increasingly used to make consequential decisions. In recent  
2 years, they have begun to advise employers on which job applicants to interview,  
3 judges on which defendants to grant bail, lenders on which homeowners to give  
4 loans, and more. In such settings, different data-driven rules result in different  
5 decisions. The problem is, for every data-driven rule, there are exceptions. While  
6 a data-driven rule may be appropriate for some, it may not be appropriate for all.  
7 In this piece, we argue that existing frameworks do not fully encompass this view.  
8 As a result, individuals are often, through no fault of their own, made to bear the  
9 burden of being data-driven exceptions. We discuss how data-driven exceptions  
10 arise and provide a framework for understanding how we can relieve the burden on  
11 data-driven exceptions. Our framework requires balancing three considerations:  
12 individualization, uncertainty, and harm. Importantly, no single consideration  
13 trumps the rest. We emphasize the importance of uncertainty, advocating that  
14 decision-makers should utilize data-driven recommendations only if the levels  
15 of individualization and certainty are high enough to justify the potential harm  
16 resulting from those recommendations. We argue that data-driven decision-makers  
17 have a duty to consider the three components of our framework before making a  
18 decision, and connect these three components to existing methods.

## 19 1 Introduction

20 We make sense of our world through rules. But, to every rule, there are exceptions. Although  
21 exceptions are by definition uncommon, they often carry significance disproportionate to their  
22 numbers. Exceptions not only improve our understanding of the rules, but they also help us develop  
23 better ones. No matter how good the rule, it cannot work for everyone, begging the question: *What*  
24 *happens to the individuals for which a decision rule is unfit: the exceptions?*

25 In some cases, nothing. We accept that rules and generalizations are, on occasion, tolerable and even  
26 necessary (Lippert-Rasmussen, 2011). Indeed, the law allows landlords to put no-pet clauses in rental  
27 agreements (a rule based on the generalization that renters with pets cause more damage to homes  
28 than renters without pets) and airlines to remove passengers for safety reasons (a policy that relies on  
29 judgments about actions that a passenger could but has not yet committed).

30 In other cases—typically, when the risk of harm is high—the state steps in to shield individuals from  
31 the adverse effects that can follow from the over-application of rules. Consider sentencing decisions.  
32 For many crimes, there are mandatory minimum sentences: a set of standardized rules that prescribe  
33 the minimum sentence a defendant must serve for a crime, if convicted. These rules arose in the  
34 U.S. as a way to “make sentencing procedures fairer and sentencing outcomes more predictable and  
35 consistent” (Travis et al., 2014). Importantly, mandatory minimum sentences were also used in capital  
36 cases, i.e., cases in which the crime is punishable by death. In 1976, however, the U.S. Supreme

37 Court ruled in *Woodson v. North Carolina* that mandatory minimum sentences should *not* be applied  
38 to capital cases. The Court wrote that there must be “consideration of the character and record of the  
39 individual offender and the circumstances of the particular offense” before imposing a sentence as  
40 serious and irrevocable as the death penalty (U.S. Supreme Court, 1976). In other words, the Court  
41 decided that, when it comes to the death penalty, rules that regularly yield exceptions—in this setting,  
42 defendants on which the rule, but not the presiding judge, would impose the death penalty—are  
43 unacceptable. The courts responded by giving greater discretion to judges.

44 In this piece, we turn our attention to *data-driven rules*. By “data-driven rules,” we refer to the  
45 decision rules behind data-driven decision aids. For example, data-driven decision aids in lending  
46 advise lenders on whether or not to grant a loan. Generally speaking, these aids produce a score  
47 for each applicant that predicts the likelihood that the applicant, if approved, would repay the loan.  
48 Different rules generate different scores. While one rule may give higher scores to applicants with  
49 families, another may not. One rule may use the applicant’s zip code as an input while another may  
50 not. As such, an applicant may be approved for a loan under some rules but not others.

51 As many scholars have acknowledged, there is a gap in the governance of data-driven decisions  
52 because individuals who are subject to data-driven decisions are not always protected by a legal  
53 system that has been built around human decisions (Citron, 2007; Wachter and Mittelstadt, 2019;  
54 Kaminski and Urban, 2021). In this piece, we focus on what happens to individuals for which a  
55 data-driven rule is unfit. We find that, under existing frameworks, *individuals on which a data-driven*  
56 *rule fails are, through no fault of their own, made to bear the burden of those mistakes*. There are two  
57 characteristics of data-driven rules that exacerbate this problem. First, data-driven rules are highly  
58 non-intuitive. Second, they are frequently updated, oftentimes without any visible changes to those  
59 that employ them.

60 This setting stands in stark contrast to how we approach rules in our legal system. As opposed to  
61 standards, rules provide a clear delineation between behaviors that are legal and those that are illegal.  
62 Rules are furthermore public and relatively static. That is, whether or not individuals agree with rules  
63 that appear in the law, they are aware of the consequences of their actions, and they do not expect that  
64 these rules change without their knowledge. Data-driven decision-subjects, on the other hand, are  
65 often left in the dark about why they receive the decision that they did, and we show that they face an  
66 disproportionately high barrier to proving that the data-driven rule is unfit for them.

67 In this piece, we unpack this phenomenon in detail, and propose a framework that remedy the problem  
68 of data-driven exceptions.

## 69 **2 Background**

70 Bringing attention to exceptions that may be neglected by systems that work well for the majority has  
71 philosophical and legal grounding. One influential concept that emphasizes the importance of giving  
72 appropriate consideration and respect to each individual is *dignity*.

73 Dignity is a concept that appears in international human rights law and domestic constitutions  
74 (O’Mahony, 2012). Despite being widely acknowledged as a “foundational principle,” its meaning  
75 and consequent role in law remain unclear (O’Mahony, 2012; Rao, 2011; Glensy, 2011). It has been  
76 used—in different and, at times, conflicting ways—to justify the right to free speech (U.S. Supreme  
77 Court, 1971), a gay couple’s right to marry (Supreme Court of California, 2008), a woman’s freedom  
78 to choose an abortion (U.S. Supreme Court, 1992), and more. Its flexible meaning allows it to serve  
79 as a unifying theoretical basis for human rights and is part of the reason it appears in the Universal  
80 Declaration of Human Rights, which states that “all human beings are born free and equal in dignity  
81 and rights” (United Nations General Assembly, 1948). Although there are multiple notions of dignity,  
82 we focus on two.

83 The first is the notion of *inherent dignity*, as popularized by Kant, who states that all humans possess  
84 “a dignity (an absolute inner worth) by which he exacts respect for himself from all other rational  
85 beings in the world” and that this dignity cannot be substituted, exchanged, gained, or lost (Kant,  
86 2017). Inherent dignity is based on the belief that, by virtue of being human, individuals must be  
87 afforded a “necessary respect” by others and the state (Gewirth, 1992). Kant (1967) also believed that  
88 individual autonomy and self-determination are special to humans and therefore intrinsically tied to  
89 dignity. In practice, inherent dignity is associated with negative liberty—a freedom from interference

90 by the state that is rooted in the idea that a “person’s dignity is best respected or enabled when he can  
91 pursue his own ends in his own way” (Rao, 2011).

92 The second notion of dignity relevant to this piece is *dignity as recognition*, which requires that  
93 there be “esteem and respect for the particularity of each individual” (Rao, 2011). It demands that  
94 an individual’s uniqueness is recognized and respected. Recall that inherent dignity is rooted in the  
95 idea that all individuals possess an inner worth that is deserving of respect regardless of whether  
96 their dignity is recognized. By contrast, under the concept of recognition dignity, an individual  
97 “can have dignity and a sense of self only through recognition by the broader society” (Rao, 2011).  
98 That all individuals inherently possess dignity is a “presumption of human equality” (Rao, 2011).  
99 On the other hand, dignity as recognition requires “treatment that *expresses* the equal worth of all  
100 individuals and their life choices” despite their differences (Rao, 2011). Rather than freedom from  
101 interference, recognition dignity is a positive concept in that the state must protect recognition dignity  
102 by enforcing respect between citizens and designing policies that actively acknowledge the equal  
103 worth of each individual (or group) in their uniqueness (Rao, 2011). In the past, recognition dignity  
104 has been invoked in claims against defamation and hate speech as well as the right to develop one’s  
105 personality (Post, 1986; Supreme Court of Canada, 1990; Federal Law Gazette, 2020).

106 The respect for an individual’s uniqueness that is demanded by recognition dignity is closely related  
107 to our work. In highlighting how the reliance of data-driven decisions on rules can inflict harm  
108 on exceptions, our work emphasizes “a basic respect for individual human dignity in a political  
109 system that otherwise allocates costs and benefits on the basis of majority rule” (Paradis, 2015).  
110 In this way, it can be viewed as a mechanism for protecting the recognition dignity of individuals  
111 in high-stakes, data-driven decision contexts. Recognizing the dignity of decision-subjects does  
112 not require that decisions always tip in their favor. It simply requires a respect for dignity—an  
113 acknowledgment that when a decision can inflict significant harm on the subject, the decision should  
114 be based on a “respectful deliberation” that balances the subject’s unique circumstances alongside  
115 other considerations (Harel, 2014).

## 116 2.1 A Note on Related Works

117 There are many existing works on data-driven technologies and their pitfalls. These works have  
118 covered enormous ground, highlighting issues that arise during the application of data-driven tech-  
119 nologies and gaps in their governance. We build on this literature, but there are four factors that  
120 together make this piece distinct.

121 First, many works (such as those examining disparate impact (Barocas and Selbst, 2016)) focus on  
122 group-based outcomes, e.g., discrimination based on a protected attribute. In contrast, our work  
123 examines data-driven decisions through the lens of individual outcomes rather than group-based ones.  
124 In particular, we discuss how one can determine if a data-driven rule is appropriate for a specific  
125 decision-subject, similarly to Lippert-Rasmussen (2011).

126 Second, most existing works propose to improve outcomes by requiring that data-driven tools be  
127 “fair,” “accurate,” and “reliable” (Citron and Pasquale, 2014; Wachter and Mittelstadt, 2019). We find  
128 that such criteria are important but do not capture the full picture when evaluating the suitability of a  
129 data-driven tool for a *specific* decision context. Accuracy, for instance, is an average notion—high  
130 accuracy only implies good performance in an average sense. Similarly, reliability implies good  
131 performance in a repeated sense—that, if run many times, an algorithm would consistently perform  
132 well. In this piece, we offer an additional consideration. In addition to fairness, accuracy, and  
133 reliability, there is another desideratum: a decision-maker should not presume that the data-driven  
134 rule is suitable for an arbitrary decision-subject, particularly when the stakes are high. Rather, the  
135 decision-maker should only apply a data-driven rule if they are sufficiently confident (as measured  
136 against the risk of harm) that it is indeed suitable, as detailed in Section 3.

137 Third, there are several works that examine whether data-driven rules should be sufficiently in-  
138 dividualized in order to be applied. That is, they investigate how individualization addresses the  
139 problem of statistical discrimination (Lippert-Rasmussen, 2011; Wachter et al., 2021). In this piece,  
140 we build on this discussion and argue that individualization is one, but not the only, component of  
141 evaluating a data-driven rule’s suitability. We maintain that one must also consider the data-driven  
142 rule’s *uncertainty*, a concept that is often overlooked but is core to the our work.

143 Lastly, our hope is to provide a framework that can serve as a common language with which to  
144 discuss data-driven exceptions across disciplines. To this end, we also consider the technical aspects  
145 of data-driven exceptions, including their origins (showing that data-driven exceptions arise in more  
146 ways than existing works typically consider, making the problem less straightforward than commonly  
147 assumed) and the technical viability of the proposed solutions. We pay particular attention to the  
148 latter. For example, although open-sourcing data-driven tools may be useful, it is infeasible in many  
149 cases (e.g., due to trade secret law or that open-sourcing introduces vulnerabilities to adversarial  
150 attacks). To ensure technical viability, we distill the our framework down to three concepts, described  
151 in Section 3, that are also meaningful in machine learning.

### 152 **3 Proposed Framework**

153 In this piece, we argue that, a decision-maker cannot presume that a data-driven rule is suitable for  
154 a given decision-subject—they must be sufficiently confident (relative to the risk of harm) that the  
155 individual is not an exception. In other words, a data-driven decision-maker—whether a machine  
156 or machine-aided human—must make a decision that inflicts harm only if they have applied due  
157 care and due diligence in determining whether the data-driven rule is fit for the decision-subject in  
158 question. The greater the risk of harm, the higher the bar.

159 Society has, for the most part, developed standards for assessing whether a human has applied due  
160 care and due diligence in decision-making (cf. the right to individualized sentencing (Berry III, 2019)).  
161 After all, the law has been honed to work for human-driven decisions. How one would operationalize  
162 this concept in the data-driven context is, however, unclear. In this piece, we propose that adapting this  
163 requirement for data-driven decisions can be achieved by considering three factors: *individualization*,  
164 *harm*, and *uncertainty*. Via these three components, we provide a concrete framework through which  
165 a decision-maker can determine when a data-driven rule is appropriate or a decision-subject can  
166 determine whether to contest a data-driven decision. Importantly, no components on its own is  
167 sufficient (or necessarily even desirable). For instance, as we unpack below, individualization can  
168 often give rise to undesirable effects.

#### 169 **3.1 Individualization: Moving from the Aggregate to the Individual**

170 For many, the natural first step to designing a data-driven rule that surpasses the appropriate levels of  
171 care and diligence in ruling out an exception is individualization: the process of tailoring a rule to  
172 the specific circumstances under consideration. In short, individualization shifts attention from the  
173 aggregate to the individual. The more individualized a rule, the more suitable it is for a particular  
174 decision-subject. For example, one way to make a data-driven rule more individualized is to add  
175 features, or inputs, to the model. A data-driven rule that uses an applicant’s age, home address, and  
176 occupation in order to decide whether to grant a loan is therefore more individualized than one that  
177 uses only their age and home address.

178 Individualization is an information concept in that it requires a decision-maker to consider the totality  
179 of an individual’s circumstances rather than make judgments based on a limited set of information.  
180 In other words, to individualize a rule is to give it additional (relevant) information. The desire  
181 for individualized decisions—the first component of our framework—is not new. Indeed, Lippert-  
182 Rasmussen (2011) discusses the right to be treated as an individual as a proposal for reducing  
183 statistical discrimination (treating an individual as if they were the statistical average of similar  
184 individuals). The push for individualization is based on the logic that, the more individualized an  
185 assessment, the less likely it is to have made broad-strokes generalizations and, as a result, to yield  
186 exceptions.

187 Individualization is a particularly useful concept because it appears in both legal texts (cf. the right to  
188 individualized sentencing (Berry III, 2019; Jorgensen, 2021)) as well as technical ones. As such, a law  
189 requiring individualization in data-driven rules would pave a clear path for computer scientists. Indeed,  
190 much of machine learning echoes the belief that, with enough information and enough historical  
191 data, a data-driven rule can predict the target outcome with perfect accuracy. Individualization has  
192 become so central to machine learning that data-driven rules are often justified based on their level of  
193 individualization. Most theorems in machine learning, for instance, follow the template: “As  $N$  goes  
194 to infinity, the error goes to 0” (occasionally accompanied by a “with high probability”), where  $N$   
195 quantifies the amount of information.

196 Perfect individualization, however, is difficult to implement. In practice, current methods are incapable  
197 of individualizing in ways that humans do naturally. Humans, for example, are generally flexible  
198 enough to update their decisions to incorporate additional pieces of information. Although a judge  
199 may initially receive certain information about a defendant, they can update their belief when given  
200 novel information (e.g., that the defendant volunteers or has dependents). Humans rely on this  
201 unique ability to holistically examine an individual’s circumstance in order to produce individualized  
202 decisions. In contrast, most (if not all) data-driven rules have fixed inputs and cannot incorporate  
203 features that are not present in the training data.

204 So, perhaps perfect individualization is not possible, but is individualizing the rule as much as possible  
205 (albeit imperfectly) all that is required to ensure that the rule is fit for use? Stated differently, suppose  
206 that a data-driven rule were perfectly individualized—that is, it incorporates all relevant information.  
207 Would such a fully individualized data-driven rule be enough?

### 208 3.2 Individualization is Not Enough: Uncertainty Also Matters

209 No—individualization is not the only pertinent factor. There are two additional components: uncer-  
210 tainty and harm, and we focus on the former in this section. The takeaway is that while individualizing  
211 a data-driven rule takes an important step toward ensuring that it does not neglect relevant information,  
212 no amount of individualization can remove all the uncertainty in a data-driven rule, and the amount  
213 of uncertainty matters when the risk of harm is high.

214 Recall that individualization is an information concept: it relies on the belief that, holding everything  
215 else equal, adding information improves a data-driven rule. Conveniently, this reasoning also underlies  
216 machine learning, which is founded on the idea that data is king (i.e., that with enough information, a  
217 data-driven rule can perform perfectly). In reality, however, even the best data-driven models make  
218 mistakes, often because some predictions are inherently impossible to get right every time. In fact,  
219 there are very few (if any) meaningful settings in which a perfect rule exists, and the main barrier is  
220 uncertainty.

221 To illustrate the limitations of individualization, consider the following two types of uncertainty  
222 (Kendall and Gal, 2017a):

- 223 1. *Epistemic uncertainty* is systematic or reducible uncertainty that arises from lack of knowledge.  
224 For example, a prediction of tomorrow’s temperature that is based on past years’ temperatures at  
225 this time of year has greater epistemic uncertainty than the prediction of tomorrow’s temperature  
226 based on past years’ temperature at this time of year *and* today’s temperature.
- 227 2. *Aleatoric uncertainty* is statistical or irreducible uncertainty that arises from the inherent ran-  
228 domness or “unknowability” of an event. At the time of prediction, no information exists that  
229 can reduce this type of uncertainty. For example, the randomness in the wind patterns that may  
230 occur between today and tomorrow prevents a temperature prediction that is made today from  
231 being perfectly certain about tomorrow’s temperature, and this randomness can be attributed to  
232 aleatoric uncertainty.

233 Through these two types of uncertainty, it becomes clear that while individualization may reduce  
234 epistemic uncertainty, it cannot reduce aleatoric uncertainty. In some cases, individualization does  
235 not even reduce epistemic uncertainty. Consider the following examples.

#### 236 **Example 1 (Individualization increases granularity at the risk of increasing uncertainty)**

237 *Consider a data scientist who wishes to increase the individualization of a data-driven rule used*  
238 *in healthcare. To do so, the data scientist adds features to the rule’s input. Instead of taking in a*  
239 *patient’s current age, height, and weight as inputs, the data scientist modifies the rule to also accept*  
240 *the patient’s history of heights and weights at every year of their life.*

241 *Suppose the data scientist uses a nearest-neighbors-style algorithm—an approach that makes a*  
242 *prediction for patient X based on previous (exemplar) patients who have similar attributes as X. Then,*  
243 *the more refined the features, the fewer exemplar patients for X exist. In other words, individualization*  
244 *reduces the amount of evidence that the nearest-neighbors rule can use to generate its assessment.*  
245 *As such, while the data scientist reduces epistemic uncertainty in one way, they increase it in another.*

246 **Example 2 (The unknowability of unobserved outcomes)** Consider a data-driven decision aid  
247 for college admissions—specifically, one that predicts how well a student will perform if admitted.  
248 Beyond random noise, there are multiple ways that aleatoric uncertainty arises.

249 For one, even if the student is similar to previous students for which there is data, one could argue that  
250 a student’s performance is not predetermined, i.e., that they have the ability to perform differently from  
251 past individuals. That each student possesses their own potential for success—that they have their  
252 own autonomy—means that no amount of individualization can predict performance with certainty.  
253 Indeed, believing that a data-driven rule carries no uncertainty holds students responsible for the  
254 performance of previous students (namely, students in the training data). While individualization  
255 can improve a data-driven prediction, it continues to hold the decision-subject responsible for the  
256 performance of previous—albeit, increasingly similar—students, and there is always uncertainty  
257 associated with the decision-subject’s own potential for success.

258 For another (and perhaps more concretely), there is also omission bias. The training data only  
259 captures the performance of students who were admitted, which implies that the performance of  
260 a student who was not admitted is unknowable (Kleinberg et al., 2017). Perhaps a student who is  
261 similar to the decision-subject but was not admitted would have performed very well.

262 Lastly, even if a decision-maker has perfect knowledge of previous students’ outcomes, any decision  
263 that is made now can only use information obtained up until this moment. There are, however,  
264 countless factors (or, in the language of causal inference, “interventions”) that could influence a  
265 student’s performance between the time of acceptance and graduation, such as whether they receive  
266 tutoring, who they befriend, and whether they take a part-time job. The only way that an assessment  
267 can be perfect and rid of uncertainty is for the target outcome itself to be an input to the assessment,  
268 but this logic is circular. If one could measure the target outcome, one would not need to infer it.

### 269 3.3 The Importance of Uncertainty: Weighing the Risk of Harm

270 In short, individualization can, at best, remove epistemic uncertainty, but no amount of individualiza-  
271 tion can remove aleatoric uncertainty. Perhaps one of the best ways to summarize this argument is via  
272 *computational irreducibility* (Wolfram, 2002). The reasoning behind this concept goes: a computer is  
273 one of many components in our world. Therefore, the complexity of a computer must be strictly lower  
274 than the complexity of the world. It follows from this logic a computer cannot predict any arbitrary  
275 outcome of interest  $Z$  (even if it was given all the historical data in the world and continually fed new  
276 data) because the complexity of the process that produces  $Z$  may be higher than the computational  
277 capacity of the computer.

278 That is not to say that data scientists should throw up their hands and give up. Indeed, computational  
279 irreducibility does not imply that every prediction task is hopeless. Rather, it says that uncertainty is  
280 inevitable when predicting a *complex* target outcome. However, eliminating uncertainty is besides the  
281 point. It is unreasonable to ask for a perfect data-driven rule that makes no mistakes. Instead, we ask  
282 that the level of uncertainty be *balanced against the risk of harm*.

283 More precisely, suppose that one of the decision outcomes would inflict significant harm. Then, no  
284 matter how individualized a decision rule may be, the decision to inflict harm should only follow if  
285 the level of certainty is high enough. If, on the other hand, the level of uncertainty (epistemic and  
286 aleatoric) is too high, then the decision-maker should err on the side of caution (less harm).

287 As an extreme example, suppose a decision-maker is presented with a newborn and must decide  
288 whether to confine them for the rest of their lives based on an evaluation of whether they will commit  
289 murder during their lifetime. The decision is made at the time of birth, so the only information that is  
290 available must also be available at the time of birth. A rule could be perfectly individualized (based  
291 on the information at the time of birth), but most would agree that there are so many unknowable  
292 factors that could contribute to the newborn’s future actions that no amount of individualization  
293 would justify inflicting a harm as high as confining a newborn for life.

294 A decision outcome’s risk of harm therefore determines the amount of individualization and certainty  
295 necessary to utilize a data-driven rule whose recommendation inflicts harm. Some decisions might  
296 carry a risk of harm so low the level of individualization and certainty needed to justify the use of a  
297 data-driven rule is accordingly low. It is natural to then ask: How should harm be measured? While  
298 providing an explicit framework for quantifying harm is out of the scope of this piece, we note that

299 prior works have laid out a path for doing so, including Wachter and Mittelstadt (2019)’s work on  
300 the right to reasonable inferences (in which they discuss the determination of “high-risk inferences”)  
301 and Kaminski and Urban (2021)’s right to contest AI (in which they characterize risk of harm in  
302 terms of “significant effects”). The European Union’s Artificial Intelligence Act also provides a “risk  
303 methodology” for categorizing high-risk decision contexts (European Union, 2022).

### 304 3.4 Putting it All Together

305 Our framework requires that the decision-maker not necessarily presume that a data-driven rule is  
306 suitable for a decision-subject, particularly when the risk of harm is high. Rather, we require that  
307 the decision-maker inflict harm only if they have applied due care and due diligence in determining  
308 whether the data-driven rule is fit for use on the individual in question. It seeks to prevent data-driven  
309 decisions from inflicting irreparable and repeated harm on individuals who, through no fault of their  
310 own, are exceptions to a data-driven rule. Our work emphasizes that data-driven rules cannot be  
311 applied blanketly. While a data-driven rule may be appropriate for some individuals, it may not be  
312 appropriate for all. In particular, when a decision may inflict significant harm on the decision-subject,  
313 examining whether or not a decision-maker is justified in using the data-driven recommendation  
314 becomes pertinent. In this way, our framework in keeping with existing concepts (originally intended  
315 for human decision-makers), including the right to dignity and the right to individualized sentencing.

316 Importantly, our framework does not imply that data-driven rules should be dropped altogether, nor  
317 does it suggest that they be used in every case. It does not even suggest that there is a clear line  
318 between the types of decisions in which data-driven rules are appropriate (e.g., that data-driven  
319 decision aids should be used in lending but not sentencing). Rather, it argues that there are some  
320 contexts in which the stakes are so high that each decision-subject deserves appropriate consideration  
321 of whether the data-driven rule is fit for them. In the same way that certain information is discarded as  
322 irrelevant (e.g., a college admissions board may discard a student’s sophomore Fall grades if a family  
323 tragedy occurred that semester), a data-driven recommendation may need to be discarded. While  
324 useful, this analogy does not carry over perfectly because it is unclear when to discard a data-driven  
325 rule. Data-driven rules behave quite differently from human ones—for instance, the “intent” and  
326 “reasoning” behind a data-driven recommendation are often inscrutable.

327 In this piece, we find that honoring an individual’s dignity requires the consideration of three factors:  
328 individualization, uncertainty, and harm. Crucially, these three factors are not only interpretable  
329 to lawmakers, but also meaningful concepts in machine learning. They therefore provide a clear  
330 language with which to assess data-driven decisions.

331 More precisely, we require that the decision-maker first evaluate the level of harm of each decision  
332 outcome. Based on the level of harm, the decision-maker should then evaluate the data-driven rule  
333 based on two considerations: individualization and uncertainty. Individualization characterizes the  
334 suitability of a rule based on how much information it considers (e.g., whether it knows enough about  
335 the decision-subject or has enough training data that pertains to the decision-subject). The level of  
336 uncertainty can be divided into two types: epistemic and aleatoric. The former captures uncertainty  
337 due to lack of information, and the latter captures the inherent unknowability of a prediction task.  
338 We require that the decision-maker utilize the data-driven recommendation only if the levels of  
339 individualization and certainty are high enough to justify the level of harm that would result from that  
340 recommendation.

## 341 4 Operationalizing the Framework

342 In this section, we examine how this framework could be operationalized. We consider what it does  
343 (and does not) mean to invoke the framework as well as *ex ante* and *ex post* measures.

### 344 4.1 Invoking the Framework

345 Does invoking the framework mean proving that the data-driven rule made a mistake? Or is it that a  
346 decision-subject who does not like their decision outcome can always claim to be an exception, thus  
347 nullifying any data-driven rule in high-risk settings?

348 Our framework says neither. Invoking the framework is *not* equivalent to proving that the data-driven  
349 rule made a mistake. For one, the outcome of interest is not always observable. In many cases, it  
350 is impossible to determine whether a mistake was made (e.g., a judge can never know whether a  
351 defendant who is denied parole would have reoffended *if* they had been granted parole instead).

352 Consider the following (simplified) example. Suppose a data-driven rule delivers random recommen-  
353 dations. For instance, suppose that it simply flips a coin each time it is asked for a recommendation.  
354 Even this random rule is bound to be correct for some individuals. However, whether this rule  
355 happens to be correct is besides the point. If the decision’s risk of harm is high (e.g., a sentencing  
356 decision), such a rule should not be applied regardless of whether or not it turns out that, down the  
357 line, the random flip happens to correctly predict the outcome. It is simply not suitable for a high-risk  
358 setting. This evaluation of a data-driven rule’s suitability is what underlies our framework. Namely,  
359 the data-driven rule should only be applied if deemed suitable for the specific decision-subject, where  
360 the level of consideration must be fitting to the risk of harm (for which we provide a framework in  
361 Section 3).

362 Importantly (and in answer to the second question above), our work does not imply that *every*  
363 individual is an exception. That is, a decision-subject who does not like their data-driven decision  
364 outcome cannot simply reverse the decision using our framework. In fact, *a data-driven rule can still*  
365 *satisfy our framework even if it makes mistakes*. It is indeed unreasonable to expect a data-driven  
366 rule to never make mistakes—a decision-subject can, at best, hope that a data-driven decision-maker  
367 ensures that a data-driven recommendation is only used if it is deemed fit for the given context. Our  
368 framework captures this principle. It does not require that a data-driven rule is perfect but that it is  
369 appropriately applied.

370 In this way, our framework is not simply a matter of mistakes. It can be violated even when a mistake  
371 has not been made (or cannot be verified). At the same time, our framework is not necessarily violated  
372 when a mistake is made. Crucially, while a data-driven rule’s accuracy—which many believe can be  
373 used to evaluate a data-driven rule’s suitability (Supreme Court of Wisconsin, 2016)—is an important  
374 performance metric, it is another way of measuring mistakes. Therefore, accuracy alone cannot fully  
375 capture the suitability of a data-driven rule, as detailed in Section 3.

## 376 4.2 *Ex ante* Justification

377 Our framework would require an *ex ante* justification that a data-driven decision appropriately  
378 considers the three components of our work—harm, individualization, and uncertainty—before  
379 such a data-driven decision is applied. Specifically, the data-driven assessment must (1) evaluate  
380 the potential harm that the decision could inflict; (2) justify the rule on the basis of its level of  
381 individualization; and (3) demonstrate that, given the level of harm and individualization, the  
382 rule appropriately and meaningfully incorporates uncertainty into its decision *or* appropriately and  
383 meaningfully communicates it to the final decision-maker.

384 In order to evaluate a decision’s potential harm, one can use the standard of “significant effects” in  
385 Article 22(1) of the General Data Protection Regulation (GDPR) (European Commission, 2016), as  
386 studied by Kaminski and Urban (2021). Wachter and Mittelstadt (2019) present a similar framework  
387 in their discussion of “high-risk inferences” with respect to the right to reasonable inferences. One  
388 can also turn to the European Union (2022)’s Artificial Intelligence Act, which provides a “risk  
389 methodology” for evaluating and categorizing high-risk decision contexts.

390 If the risk of harm is high enough, the next step is to characterize a data-driven rule’s level of  
391 individualization (which may vary across decision-subjects), as given by (2). Characterizing a rule’s  
392 level of individualization can be done in multiple ways. For example, one could require that a  
393 data-driven decision aid report its input variables, which reflect the data-driven rule’s granularity.  
394 One could also require that the data-driven rule be evaluated on performance metrics that are more  
395 fine-grained than accuracy, such as calibration or multicalibration (Hébert-Johnson et al., 2018).

396 Lastly, (3) is a final step that combines the insights of (1) and (2). Specifically, (3) determines  
397 whether, given the potential harm assessed in (1) and the level of individualization found in (2), the  
398 final decision appropriately and meaningfully incorporates uncertainty. If the final decision maker is  
399 the *algorithm*, one must demonstrate that the assessment appropriately and meaningfully considers  
400 uncertainty. If the final decision-maker is *human*, then the data-driven assessment must appropriately  
401 and meaningfully communicate uncertainty to them. Incorporating uncertainty is necessary, as it



402 synthesizes the assessments of harm and individualization. For instance, if a decision carries a risk of  
403 significant harm, then the level of individualization and the accompanying certainty may not be high  
404 enough to justify inflicting harm. It may even be the case that, in certain contexts, no matter how  
405 individualized the assessment, there is too much uncertainty to justify inflicting harm while, in others,  
406 the risk of harm is so low that a high level of uncertainty is acceptable. Meaningfully incorporating  
407 or communicating uncertainty for (3) is an active area of research in human-computer interaction  
408 (Hullman, 2016; Hofman et al., 2020). To communicate uncertainty meaningfully, the assessment  
409 could report on different types of uncertainty, similarly to how existing works distinguish between  
410 epistemic and aleatoric uncertainty (Kendall and Gal, 2017b), as discussed in Section 3.

### 411 4.3 *Ex post* Contestation

412 It is important that decision-subjects be able to contest a data-driven decision *ex post*. As explored by  
413 Kaminski and Urban (2021), contestation is an accountability mechanism that enhances the legitimacy  
414 of data-driven assessments as well as builds the public’s trust in them.

415 As a possible template, one could turn to the procedure for contesting on the basis of Title VII of  
416 the U.S. Civil Rights Act’s notion of disparate impact (Barocas and Selbst, 2016). Specifically, Title  
417 VII prohibits employment discrimination due to the individual’s race, color, religion, sex, or national  
418 origin. A plaintiff—an individual who believes that their (potential) employer violated Title VII—can  
419 sue the employer by providing evidence of what is known as “disparate impact.”

420 In disparate impact cases, a plaintiff must first establish that an employment practice negatively  
421 impacts a class of individuals protected by Title VII compared to its impact on individuals outside  
422 the protected class. Even if disparate impact is established, however, it can be countered if the  
423 defendant—or employer—successfully shows that the employment practice is rooted in “business  
424 necessity.” The defense of “business necessity” can further be refuted if the plaintiff provides a  
425 compelling alternate employment practice that would mitigate disparate impact without violating  
426 business necessity. Contestation on the basis of our framework could mirror this three-stage procedure,  
427 as follows. First, the plaintiff must establish that (1), (2), and/or (3) from Section 4.2 has been violated  
428 by the data-driven decision. If the plaintiff is successful, the defendant can counter by showing  
429 that the data-driven decision could not have been changed without demanding significant resources  
430 or inflicting disproportionate harm on other parties. Finally, if the defendant is successful in this  
431 second stage, the plaintiff can refute the defendant’s justification by providing an alternate procedure  
432 that improves upon the assessment with respect to (1)-(3) and does not demand excessive resources  
433 or inflict disproportionate harm on other parties. This procedure is one among many contestation  
434 mechanisms, as surveyed by Kaminski and Urban (2021).

## 435 5 Conclusion

436 It is widely acknowledged that the governance of data-driven decisions requires new concepts and  
437 tools. In this work, we argue that decision-subjects are often, through no fault of their own, made to  
438 bear the burden of imperfect data-driven rules. While we cannot data-driven rules are perfect, there  
439 are several characteristics of data-driven decision-making that require special treatment. In particular,  
440 data-driven rules are not only unintuitive, but also frequently updated. As such, a persistent problem  
441 within data-driven decision-making is that it will be difficult to detect when a data-driven rule makes  
442 mistakes, imposing a high burden on decision-subjects who are made to bear the cost of being one of  
443 those individuals, i.e., of being exceptions.

444 In this piece, we argue that the presumption should not be that a data-driven rule—even one that  
445 has high accuracy—is suitable for an arbitrary decision-subject of interest. Rather, a decision-maker  
446 should only apply a data-driven rule if they have applied due care and due diligence (relative to  
447 the risk of harm) in excluding the possibility that the decision-subject is an exception to the given  
448 data-driven rule. In some cases, the risk of harm may be so low that only cursory consideration is  
449 required. In others, the risk of harm may be so high that a decision-maker must be convinced that  
450 the data-driven rule works well on the *specific* decision-subject of interest before applying it. We  
451 provide a three-part framework—that requires balances individualization, harm, and uncertainty—for  
452 determining whether a data-driven decision is fit for the decision-subject of interest.

453 **References**

- 454 Barocas, S. and A. D. Selbst (2016). Big Data’s Disparate Impact. *California Law Review* 104, 671.
- 455 Berry III, W. W. (2019). Individualized Sentencing. *Washington & Lee Law Review* 76, 13.
- 456 Citron, D. K. (2007). Technological Due Process. *Washington University Law Review* 85, 1249.
- 457 Citron, D. K. and F. Pasquale (2014). The Scored Society: Due Process for Automated Predictions.  
458 *Washington Law Review* 89, 1.
- 459 European Commission (2016). General Data Protection Regulation.
- 460 European Union (2022). The Artificial Intelligence Act.
- 461 Federal Law Gazette (2020). *Grundgesetz Für die Bundesrepublik Deutschland [Basic Law for the*  
462 *Federal Republic of Germany]*.
- 463 Gewirth, A. (1992). Chapter 1: Human Dignity as the Basis of Rights. In M. J. Meyer and W. A.  
464 Parent (Eds.), *The Constitution of Rights*, pp. 10–28. Cornell University Press.
- 465 Glensy, R. D. (2011). The Right to Dignity. *Columbia Human Rights Law Review* 43, 1–65.
- 466 Harel, A. (2014). *Why Law Matters*. Oxford Legal Philosophy. Oxford University Press.
- 467 Hébert-Johnson, U., M. Kim, O. Reingold, and G. Rothblum (2018). Multicalibration: Calibration  
468 for the (Computationally-Identifiable) Masses. In *International Conference on Machine Learning*,  
469 pp. 1939–1948. PMLR.
- 470 Hofman, J. M., D. G. Goldstein, and J. Hullman (2020). How Visualizing Inferential Uncertainty  
471 Can Mislead Readers About Treatment Effects in Scientific Results. In *Proceedings of the 2020*  
472 *Conference on Human Factors in Computing Systems (CHI)*, pp. 1–12.
- 473 Hullman, J. (2016). Why Evaluating Uncertainty Visualization is Error Prone. In *Proceedings of the*  
474 *Sixth Workshop on Beyond Time and Errors on Novel Evaluation Methods for Visualization*, pp.  
475 143–151.
- 476 Jorgensen, R. (2021). Algorithms and the Individual in Criminal Law. *Canadian Journal of*  
477 *Philosophy*, 1–17.
- 478 Kaminski, M. E. and J. M. Urban (2021). The Right to Contest AI. *Columbia Law Review* 121(7),  
479 1957–2048.
- 480 Kant, I. (1967). *Grundlegung zur Metaphysik der Sitten*. Universal-Bibliothek. Reclam.
- 481 Kant, I. (2017). *The Metaphysics of Morals*. Cambridge Texts in the History of Philosophy. Cambridge  
482 University Press.
- 483 Kendall, A. and Y. Gal (2017a). What Uncertainties Do We Need in Bayesian Deep Learning for  
484 Computer Vision? *Advances in neural information processing systems* 30.
- 485 Kendall, A. and Y. Gal (2017b). What Uncertainties Do We Need in Bayesian Deep Learning for  
486 Computer Vision? In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,  
487 and R. Garnett (Eds.), *Advances in Neural Information Processing Systems*, Volume 30. Curran  
488 Associates, Inc.
- 489 Kleinberg, J., H. Lakkaraju, J. Leskovec, J. Ludwig, and S. Mullainathan (2017). Human Decisions  
490 and Machine Predictions. *The Quarterly Journal of Economics* 133(1), 237–293.
- 491 Lippert-Rasmussen, K. (2011). “We Are All Different”: Statistical Discrimination and the Right to  
492 be Treated as an Individual. *The Journal of ethics* 15(1), 47–59.
- 493 O’Mahony, C. (2012, 03). There is no such thing as a right to dignity. *International Journal of*  
494 *Constitutional Law* 10(2), 551–574.
- 495 Paradis, M. (2015, December). Dignity through Law. The New Rambler.

- 496 Post, R. C. (1986). The Social Foundations of Defamation Law: Reputation and the Constitution.  
497 *California Law Review* 74, 691.
- 498 Rao, N. (2011). Three Concepts of Dignity in Constitutional Law. *Notre Dame Law Review* 86,  
499 1–183.
- 500 Supreme Court of California (2008). In re Marriage Cases, 43 Cal.4th 757, 76 Cal. Rptr. 3d 683, 183  
501 P.3d 384.
- 502 Supreme Court of Canada (1990). R. v. Keegstra, 3 S.C.R. 697.
- 503 Supreme Court of Wisconsin (2016). State v. Wisconsin, 881 N.W.2d 749, 754, 757.
- 504 Travis, J., B. Western, and F. S. Redburn (2014). The Growth of Incarceration in the United States:  
505 Exploring Causes and Consequences.
- 506 United Nations General Assembly (1948). Universal Declaration of Human Rights.
- 507 U.S. Supreme Court (1971). Cohen v. California, 403 u.s. 15.
- 508 U.S. Supreme Court (1976). Woodson v. North Carolina, 428 U.S. 280, 304.
- 509 U.S. Supreme Court (1992). Planned Parenthood of Southeastern PA. v. Casey, 505 U.S. 833.
- 510 Wachter, S. and B. Mittelstadt (2019). A Right to Reasonable Inferences: Re-Thinking Data Protection  
511 Law in the Age of Big Data and AI. *Columbia Business Law Review*, 494.
- 512 Wachter, S., B. Mittelstadt, and C. Russell (2021). Why fairness cannot be automated: Bridging the  
513 gap between eu non-discrimination law and ai. *Computer Law & Security Review* 41, 105567.
- 514 Wolfram, S. (2002). *A New Kind of Science*. Wolfram Media.