

STEREOTYPES

PEDRO BORDALO

KATHERINE COFFMAN

NICOLA GENNAIOLI

ANDREI SHLEIFER*

December 14, 2015

Abstract

We present a model of stereotypes based on Kahneman and Tversky’s representativeness heuristic. A decision maker assesses a target group by overweighting its representative types, which we formally define to be the types that occur more frequently in that group than in a baseline reference group. Stereotypes formed in this way contain a “kernel of truth”: they are rooted in true differences between groups. They are also context dependent: beliefs about a group depend on the characteristics of the reference group. Because stereotypes emphasize differences, they cause belief distortions, particularly when groups are similar. In line with our predictions, beliefs in the lab about abstract groups and beliefs in the field about political groups are context dependent and distorted in the direction of representative types.

*Royal Holloway University of London, Ohio State University, Università Bocconi and IGER, Harvard University. We are grateful to Nick Barberis, Roland Bénabou, Dan Benjamin, Tom Cunningham, Matthew Gentzkow, Emir Kamenica, Larry Katz, David Laibson, Sendhil Mullainathan, Josh Schwartzstein, Jesse Shapiro, Alp Simsek, Neil Thakral, and four anonymous referees for extremely helpful comments, and Jesse Graham, Jonathan Haidt, and Brian Nosek for sharing the Moral Foundations Questionnaire data. We thank the Initiative on Foundations of Human Behavior for support of this research, and Maik Wehmeyer and Laura Freitag for research assistance. Corresponding Author: Andrei Shleifer, email: ashleifer@harvard.edu, tel: 6174955046, Littauer Center, 1805 Cambridge Street, Cambridge, MA 02138. JEL: D03, D83, D84.

1 Introduction

The Oxford English Dictionary defines a stereotype as a “widely held but fixed and oversimplified image or idea of a particular type of person or thing”. Stereotypes are ubiquitous. Among other things, they cover racial groups (“Asians are good at math”), political groups (“Republicans are rich”), genders (“Women are bad at math”), demographic groups (“Florida residents are elderly”), and situations (“Tel-Aviv is dangerous”). As these and other examples illustrate, some stereotypes are roughly accurate (“the Dutch are tall”), while others much less so (“Irish are red-headed”; only 10% are). Moreover, stereotypes change: in the US, Jews were stereotyped as religious and uneducated at the beginning of the 20th century, and as high achievers at the beginning of the 21st (Madon et. al., 2001).

Social science has produced three broad approaches to stereotypes. The economic approach of Phelps (1972) and Arrow (1973) sees stereotypes as a manifestation of statistical discrimination: rational formation of beliefs about a group member in terms of the aggregate distribution of group traits. Statistical discrimination may impact actual group characteristics in equilibrium (Arrow 1973), but even so stereotypes are based on rational expectations.¹ As such, these models do not address the central problem that stereotypes are often inaccurate. The vast majority of Florida residents are not elderly, the vast majority of the Irish are not red-headed, and Tel-Aviv is really pretty safe.

The sociological approach to stereotyping pertains only to social groups. It views stereotypes as fundamentally incorrect and derogatory generalizations of group traits, reflective of the stereotyper’s underlying prejudices (Adorno et al. 1950) or other internal motivations (Schneider 2004). Social groups that have been historically mistreated, such as racial and ethnic minorities, continue to suffer through bad stereotyping, perhaps because the groups in power want to perpetuate false beliefs about them (Steele 2010, Glaeser 2005). The stereotypes against blacks are thus rooted in the history of slavery and continuing discrimination. This approach might be relevant in some important instances, but it leaves a lot out. While some stereotypes are inaccurate, many are quite fair (“Dutch are tall,” “Swedes

¹More recent work suggests that stereotypes are not necessarily self-fulfilling: assuming that freely available information is used correctly, minorities can invest in visible signals of quality that offset preconceptions (Lundberg and Startz 1983). Glover et al (2015) present evidence on self-fulfilling aspects of stereotypes.

are blond.”) Moreover, many stereotypes are flattering to the group in question rather than pejorative (“Asians are good at math”). Finally, stereotypes change, so they are at least in part responsive to reality rather than entirely rooted in the past (Madon et. al., 2001).

The third approach to stereotypes – and the one we follow – is the “social cognition approach”, rooted in social psychology (Schneider 2004). This approach gained ground in the 1980’s and views social stereotypes as special cases of cognitive schemas or theories (Schneider, Hastorf, and Ellsworth 1979). These theories are intuitive generalizations that individuals routinely use in their everyday life, and entail savings on cognitive resources. Hilton and Hippel (1996) stress that stereotypes are “mental representations of real differences between groups [. . .] allowing easier and more efficient processing of information. Stereotypes are selective, however, in that they are localized around group features that are the most distinctive, that provide the greatest differentiation between groups, and that show the least within-group variation.” A related “kernel-of-truth hypothesis” holds that stereotypes are based on some empirical reality; as such, they are useful, but may entail exaggerations (Judd and Park 1993).

We show that this approach to stereotypes is intimately related to another idea from psychology: the use of heuristics in probability judgments (Kahneman and Tversky 1972). Just as heuristics simplify the assessment of complex probabilistic hypotheses, they also simplify the representation of heterogeneous groups, sometimes causing errors in judgment. We formally explore this idea by modelling stereotype formation as a consequence of Kahneman and Tversky’s representativeness heuristic. Tversky and Kahneman (1983) write that “an attribute is representative of a class if it is very diagnostic; that is, the relative frequency of this attribute is much higher in that class than in the relevant reference class.” Following Gennaioli and Shleifer (GS 2010), we assume that a type t is representative for group G if - in line with the Tversky and Kahneman definition - it scores high on the likelihood ratio:

$$\frac{\Pr(t|G)}{\Pr(t|-G)}. \tag{1}$$

The most representative types come to mind first, and so are overweighted in judgments. Predictions about G are then made under a distorted distribution, or stereotype,

that overweights representative types. Our results obtain with minimal assumptions on such overweighting. We describe a number of weighting specifications and explore their properties.

To illustrate the logic of the model, consider the stereotype “Florida residents are elderly”. The proportion of elderly people in Florida and in the overall US population is shown below.²

| <i>age</i> | 0 – 19 | 20 – 44 | 45 – 64 | 65+ |
|------------|--------|---------|---------|-------|
| Florida | 24.0% | 31.7% | 27.0% | 17.4% |
| US | 26.9% | 33.6% | 26.4% | 13.1% |

The table shows that the age distributions in Florida and in the rest of the US are very similar. Yet, someone over 65 is highly representative of a Florida resident, because this age bracket maximizes the likelihood ratio $\Pr(t|\text{Florida})/\Pr(t|\text{US})$. When thinking about the age of Floridians, then, the “65+” type immediately comes to mind because in this age bracket Florida is most different from the rest of the US, in the precise sense of representativeness. Representativeness-based recall induces an observer to overweight the “65+” type in his assessment of the average age of Floridians.

This example also illustrates how stereotypes can be inaccurate. Indeed, and perhaps surprisingly, only about 17% of Florida residents are elderly. The largest share of Florida residents, nearly as many as in the overall US population, are in the age bracket “19-44”, which maximizes $\Pr(t|\text{Florida})$. Being elderly is not the most likely age bracket for Florida residents, but rather the age bracket that occurs with the highest relative frequency. A stereotype-based prediction that a Florida resident is elderly has very little validity.

The same logic of representativeness suggests that the reason people stereotype the Irish as red-headed is that red hair is more common among the Irish than among other groups, even though it is not that common in absolute terms. The reason people stereotype Republicans as wealthy is that the wealthy are more common among Republicans than Democrats.³ In both cases, the representation entails judgment errors: people overestimate the proportion of red-haired among the Irish, or of the wealthy among the Republicans.

We find that representativeness often generates fairly accurate stereotypes but sometimes

²Data from the 2010 US Census, see http://factfinder.census.gov/faces/tableservices/jsf/pages/productview.xhtml?pid=DEC_10_DP_DPDP1&src=pt.

³See www.nytimes.com/packages/pdf/politics/20041107_px_ELECTORATE.xls.

causes stereotypes to be inaccurate, particularly when groups have similar distributions that differ most in unlikely types. More generally, our model highlights two critical properties:

- Stereotypes amplify systematic differences between groups, even if these differences are in reality very small. When groups differ by a shift in means, stereotyping exaggerates differences in means, and when groups differ by an increase in variance, stereotyping exaggerates the differences in variances. In these cases (but not always), representativeness yields stereotypes that contain a “kernel of truth”, in the sense that they differentiate groups along existing and highly diagnostic characteristics, exactly as Hilton, Hippel and Schneider define them. Accordingly, as group traits change, stereotypes also change, consistent with Madon et al. (2001).
- Stereotypes are context dependent. The assessment of a given target group depends on the reference group to which it is compared.

In line with the social cognition approach to stereotypes, a significant body of psychological research on beliefs about gender, race, age and political groups, finds that stereotypes broadly reflect reality but also display biases (for a review, see Jussim et al. 2015). Recent work has highlighted the fact that beliefs exaggerate group differences, particularly in the context of stereotypes about age and political groups (Chan et al 2012, Westfall et al 2015). This descriptive work provides the backdrop for our own empirical investigation, which tests the properties of representativeness-based stereotypes.

We first assess the role of representativeness and context dependence in the lab. We construct a group of mundane objects, G , and present it to participants next to a comparison group, $-G$. In our baseline condition, the comparison group is chosen so that no type is particularly representative of group G . In our treatment, we change the comparison group, $-G$, while leaving the target group, G , unchanged. The new comparison group gives rise to highly representative types within G . In line with the key prediction of our model, participants in the treatment condition shift their assessment of G toward the new representative types.

We next test the model using field data. We use two data sets on political preferences, and beliefs about political preferences, in the U.S. Here, groups are political constituencies

(Democrats and Republicans) and types are their positions on an issue. We first show that beliefs depart from the truth by exaggerating (mean) differences, as per the kernel of truth logic. We then show that distortions in beliefs can be accounted for by overweighting types that are representative of each political group, in the context of the other group.

In both data sets, individual beliefs systematically depart from rational expectations and these departures can be accounted for by representativeness as we model it here. While representativeness is not the only heuristic that shapes recall (availability, driven by recency or frequency of exposure, also plays a role), it explains the fact that, in the data, stated beliefs indeed exaggerate differences among groups.

Since Kahneman and Tversky's (1972, 1973) work on heuristics and biases, several studies have formally modelled heuristics about probabilistic judgments and incorporated them into economic models. Work on the confirmation bias (Rabin and Schrag 1999) and on probabilistic extrapolation (Grether 1980, Barberis, Shleifer, and Vishny 1998, Rabin 2002, Rabin and Vayanos 2010, Benjamin, Rabin and Raymond 2011) assumes that the decision maker has an incorrect model in mind or incorrectly processes available data. Our approach is instead based on the assumption that representative information comes foremost to mind when making judgments. The mental operation that lies at the heart of our model – generating a prediction for the distribution of types in a group, based on data stored in memory – also captures base-rate neglect and confirmation bias. The underweighting or neglect of information in our model simplifies judgment problems in a way related to models of categorization (Mullainathan 2002, Fryer and Jackson 2008). In these models, however, decision makers use coarse categories organized according to likelihood, not representativeness. This approach generates imprecision but does not create a systematic bias for overestimating unlikely events, nor does it allow for context dependent beliefs. In our empirical analysis of political beliefs, we explicitly compare the predictions of representativeness-based recall to those of likelihood based models and find that the evidence supports the former.

In modeling representativeness we follow the specification of Gennaioli and Shleifer (GS 2010), but investigate a new set of questions. GS (2010) examine how representativeness distorts the assessed probabilities of alternative hypotheses, but not how the probability of a given hypothesis or group is distributed across its constituent elements. In the context

of the current setting, GS (2010) ask how imperfect recall affects the assessed probability that a randomly drawn member from a universe Ω belongs to group G . The current paper, in contrast, asks which type t we expect to draw *once we know* that we are facing group G . GS (2010) show how representativeness generates biased probabilistic assessments such as conjunction and disjunction fallacies. The current paper deals with perhaps a broader and more ubiquitous problem of stereotype formation, extensively studied by other social scientists but largely neglected by economists.

Section 2 describes our model. In Section 3 we examine the properties of stereotypes, including the forces that shape stereotype accuracy, and illustrate these properties with a number of examples. In Section 4 we bring the model to the data, by performing a lab experiment and analyzing existing surveys of political beliefs. Section 5 concludes. Appendix A contains the proofs. The Online Appendix presents a number of extensions of the model, as well as additional results for the experiments and field evidence.

2 A Model of Representativeness and Stereotypes

2.1 The Model

A decision maker (DM) faces a *prediction* problem, such as assessing the ability of a job candidate coming from a certain ethnic group, the future performance of a firm belonging to a certain sector, or future earnings based on own gender.

Formally, there is a set of types of interest T and an overall population Ω , of which group G is a subset. The set of types T can be unordered (e.g., occupations) or ordered (and typically cardinal, e.g., earnings levels). When T is ordered, we write $T = \{t_1, \dots, t_T\}$ with $t_1 < t_2 < \dots < t_T$.⁴ There is a probability or frequency distribution $\pi : T \times \Omega \rightarrow [0, 1]$, that induces a conditional distribution $\Pr(T = t | G)$ when restricted to G .⁵ In what follows, we

⁴We denote the number of types $|T|$ by T . The model applies also to cases in which types: i) are multi-dimensional, capturing a bundle of attributes such as occupation and nationality, or ii) are continuous. We consider these cases in Appendices C and D respectively. Also, G may represent any category of interest, such as the historical performance of a firm or industry, actions available to a decision maker ($T =$ set of payoffs, $G =$ occupations), or categories in the natural world ($T =$ ability to fly, $G =$ birds).

⁵In many applications each individual in Ω is characterized by a deterministic type (e.g. age, hair color, etc). As a result, $\pi(t, \omega) = 1/|\Omega|$. For instance, each Floridian has a single age type (at the finest temporal

denote by $\pi_{t,G} = \Pr(T = t | G)$ the probability of type t in group G and by π_G the vector $(\pi_{t,G})_{t \in T}$ containing the conditional distribution.

The DM's goal is to assess the distribution of the types of interest in a particular group G . While the DM has stored in memory the full distribution, he retrieves from memory a distorted version of π_G that overweights the probability of those types that are most representative of G relative to the comparison group $-G = \Omega \setminus G$. Definitions 1 and 2 formalize this representativeness-based recall, following GS (2010).

Definition 1 *The representativeness of type t for group G is defined as the likelihood ratio:*

$$R(t, G) = \frac{\pi_{t,G}}{\pi_{t,-G}}. \quad (2)$$

In line with Tversky and Kahneman (1983), a type t is representative of G if it is relatively more likely to occur in G than in $-G$. The representative age of a Floridian is 65+ because people in this age bracket are more common in Florida as compared to the rest of the US. Definition 1 implies that DMs are attuned to log differences in probabilities: representativeness depends on the percentage probability increase of a type from $-G$ to G . This captures a form of diminishing sensitivity, whereby, for a fixed probability difference, a type is more likely to be overweighted if it is unlikely in the comparison group.⁶ Statistically, representative types are also diagnostic of the target group G . Indeed, the higher is $R(t, G)$, the more confident is a Bayesian DM observing t that t belongs to G rather than to $-G$.⁷

The ease of recall of highly representative types affects judgments because more easily recalled types are overweighted. We model distorted recall as follows. Denote by $\pi_G/\pi_{-G} \equiv (\pi_{t,G}/\pi_{t,-G})_{t \in T}$ the vector of representativeness of all types in G . We then have:

resolution). When instead types are stochastic, such as when estimating future earnings of a person or a firm, each individual is described by a non-degenerate distribution.

⁶Our definition of representativeness links to Weber's law of sensory perception, see Section 2.2. It also links to our previous work on salience, in which we postulated that log differences in payoffs determine the attention to lottery payoffs, Bordalo, Gennaioli, and Shleifer (2012) and to goods' attributes (Bordalo, Gennaioli, and Shleifer 2013). Equation (2) establishes the same principle for the domain of probabilities.

⁷This insight led Tenenbaum and Griffiths (2001) to define representativeness as individuals' sense, as intuitive Bayesians, of updating in reaction to data. Their definition, like ours, is in terms of the likelihood ratio. However, Tenenbaum and Griffiths interpret representativeness as a mechanism that affects intuitive judgments of similarity, rather than beliefs (e.g. it accounts well for lab evidence where subjects are asked to rank types in terms of representativeness, or of strength of association with a group.) Accordingly, they do not consider the possibility of systematically distorted, and context dependent, beliefs.

Definition 2 *The DM attaches to each type $t \in T$ in group G a distorted probability:*

$$\pi_{t,G}^{st} = \pi_{t,G} \frac{h_t(\pi_G/\pi_{-G})}{\sum_{s \in T} \pi_{s,G} h_s(\pi_G/\pi_{-G})}, \quad (3)$$

where $h_t : R_+^T \rightarrow \mathbb{R}_+$ is a weighting function such that:

1) *Weights depend only on representativeness. Formally, $h_t = h\left(\frac{\pi_{t,G}}{\pi_{t,-G}}; \left(\frac{\pi_{s,G}}{\pi_{s,-G}}\right)_{s \in T \setminus \{t\}}\right)$ where $h : R_+ \times R_+^{T-1} \rightarrow \mathbb{R}_+$ is a function that is invariant to a permutation of the last $T - 1$ arguments.*

2) *Weighing of a type increases in own representativeness and decreases in the representativeness of other types. Formally, the function $h(\cdot)$ is weakly increasing in its first argument, and weakly decreasing in the other $T - 1$ arguments.*

We call the distribution $(\pi_{t,G}^{st})_{t \in T}$ the stereotype for G . If a type t is objectively more likely, namely $\pi_{t,G}$ is higher, then the stereotype attaches higher probability to it. By property 1), distortions are due exclusively to the fact that a type is more or less representative than the others. In particular, if all types are equally representative, the DM equally weighs all of them at $h(1)$ and holds rational expectations about G . If instead the representativeness of different types differs, property 2) implies that the stereotype ceteris paribus overweights the probability of more representative types.

Most of the results we explore in this paper hold for a general weighing function $h_t(\cdot)$. Specific functional forms capture added assumptions about the psychology of representativeness-based recall, and are useful in applications. We outline a few specifications and their properties.

- Rank-based stereotypes: the ranking of the representativeness of different types shapes distortions. Denote by $r(t) \in \{1, \dots, T\}$ the representativeness ranking of type t . When $r(t) = 1$ type t is the most representative one (potentially with ties). We can specify two ways in which a type's representativeness ranking distorts its probability.
 - Rank-based truncation: the DM only recalls the types that have representativeness ranking of at most d , namely $\{t \in T \mid r(t) \leq d\}$. Zero probability is attached

to the remaining types.⁸ Denote by $I(r(t) \leq d)$ an indicator function taking value 1 if the representativeness ranking of t is at most d . Then, the weighting function is $h_t = I(r(t) \leq d)$ so that:

$$\pi_{t,G}^{st} = \pi_{t,G} \frac{I(r(t) \leq d)}{\sum_{s \in T} \pi_{s,G} I(r(s) \leq d)},$$

which is the true conditional probability within recalled types. This assumption is used in Gennaioli and Shleifer (2010).⁹

- Rank-based discounting: The DM discounts by a constant factor $\delta \in [0, 1]$ the odds of type t relative to its immediate predecessors in the representativeness ranking. Lower δ implies stronger discounting of less representative types. Formally, the weighting function is $h_t = \delta^{r(t)}$, so that:

$$\pi_{t,G}^{st} = \pi_{t,G} \frac{\delta^{r(t)}}{\sum_{s \in T} \pi_{s,G} \delta^{r(s)}}.$$

- Representativeness based discounting: All else equal, the weight attached by the DM to type t increases continuously with its representativeness. One convenient formulation is $h_t = (\pi_{t,G}/\pi_{t,-G})^\theta$ so that:

$$\pi_{t,G}^{st} = \pi_{t,G} \frac{(\pi_{t,G}/\pi_{t,-G})^\theta}{\sum_{s \in T} \pi_{s,G} (\pi_{s,G}/\pi_{s,-G})^\theta},$$

where $\theta \geq 0$ captures the extent to which representativeness distorts beliefs. This formulation is particularly convenient when dealing with continuous distribution of the exponential or power classes.

These functional forms all embody the main idea of our model that the stereotype over-

⁸These neglected types are not viewed as impossible; they are just assigned zero probability in the DM's current thinking. This formulation allows us to model surprise and reactions to unforeseen contingencies, which have proved useful ingredients in modeling probabilistic judgments (GS 2010) as well as neglect of risk in financial crises (Gennaioli, Shleifer, and Vishny 2012).

⁹Specifically, in GS (2010) the assessed probability that a certain hypothesis G is true is equal to:

$$\Pr(G) = \frac{\sum_t \pi_{t,G} I(r(t) \leq d)}{\sum_t \pi_{t,G} I(r(t) \leq d) + \sum_t \pi_{t,\Omega/G} I(r(t) \leq d)}$$

which increases in the ratio between the total probability mass recalled for G and that recalled for $-G = \Omega \setminus G$.

weights the probability of more representative types. Rank-based truncation captures a central manifestation of limited memory: forgetting unrepresentative types. Smoother discounting (based on ranking or on representativeness) may be more appropriate when the type space is small, and smooth discounting can be more tractable in certain settings.

Section 3 characterizes the general properties of stereotypes. In particular, it shows their ability to account for social psychologists’s “kernel of truth” hypothesis under the general weighting function of Definition 2. To bring the model to the data in Section 4.2, we derive linear approximations of stereotypical beliefs by assuming that the weighting function is differentiable with respect to a type’s representativeness. This assumption excludes rank-based weighting but allows for many possibilities.

2.2 Discussion of Assumptions

Before moving to the formal analysis, we discuss some properties as well as limitations of our approach. Representativeness-based recall, the idea that individuals recall distinctive group types, can be viewed as an instance of what Kahneman and Tversky call “attribute substitution”. When dealing with the difficult question “what is the distribution of hair color among the Irish?”, people intuitively answer to the simpler question “which hair color distinguishes the Irish people?”. Critically, as discussed by Kahneman and Tversky, attribute substitution does not occur because people misunderstand the original question, or mechanically confuse the assessment of $\Pr(t|G)$ with that of $\Pr(G|t)$. Rather, it occurs because the distinctive or representative types immediately come to mind, and individuals anchor their overall probability judgment to it. As a consequence, subjects do not only make mistakes in judging the probability that a Floridian is over 65. They also give too high an answer to the question “what is the average age of a Floridian?”

One interesting question is whether the process of stereotyping we describe is optimal in some sense. We do not formally analyze this question here, but cognitive psychology offers some relevant considerations. A key operating principle of human visual perception is the highlight of contrast. An object is perceived to be brighter if set against a darker background, and vice versa. The contrast principle in visual perception has been justified as an optimal way to identify brightness, color, size, distance, in the presence of multiplicative

background noise (Kersten et al. 2004, Cunningham 2013). A similar argument can be made with respect to cognitive operations. Kahneman and Tversky invoke the contrast principle of sensory perception to justify the Prospect Theory assumption that the carriers of utility are changes with respect to a reference point. Our formulation of representativeness is related to the same idea, in the sense that individual perceptions stress differences between groups. In a noisy world in which attention is limited, this process may optimally allow for swift reactions to changes in group characteristics, even if sometimes errors are made. Exploring this idea formally is an interesting avenue for future work.

Consider now some limitations of our model. First, representativeness is not the only heuristic that shapes recall. Decision makers may for instance find it easier to recall types that are sufficiently likely. Another potentially important mechanism is availability, understood by Kahneman and Tversky (1972) as the “ease” with which information comes to mind (because of actual frequency or repetition). In Online Appendix E we present a truncation-based recall mechanism in which distortions are driven by a combination of representativeness and likelihood of types (which is equivalent to relaxing property 1 in Definition 2). This model can offer a useful starting point to capture availability as well, even though a full model of availability is beyond the scope of this paper. Even in this more general setting, the influence of representativeness on recall is the driving force of stereotypes that, in line with the social psychology perspective, are based on underlying differences among groups. As we show in Section 4, this feature is critical in accounting for the evidence.

The second set of model-related issues concerns how to specify the elements of Definition 1 in applications: group G , the type space T , and the reference group $-G$. Take the specification of the group G and of the type space T . Often, the problem itself provides a natural specification of these features. This is the case in the empirically important class of “closed end” questions, such as those used in surveys, which provide respondents with a set of alternatives, as in the data we use in Section 4. More generally, the problem solved by the decision maker – such as evaluating the resume of a job applicant coming from a certain ethnic group – primes a group, a dimension of interest, and a set of types (e.g., the applicant’s qualification or skill levels). When types have a natural order, such as income, age, or education, the granularity of T is also naturally given by the problem

(income, age, and years of schooling brackets). When the set of types is not specified by the problem, decision makers spontaneously generate one.¹⁰ It would be useful to have a model of which dimensions and types come to mind, particularly for more open ended problems. Psychologists have sought to construct a theory of natural types and dimensions (Rosch 1998). We do not make a contribution to this problem, but note that in many problems of interest in economics the dimension as well as the set of types is naturally given. Furthermore, in our model details of the type space can be important under rank-based truncation, but they matter less under smooth discounting.

Consider finally the role of the comparison group $-G$. This group captures the context in which a stereotype is formed and, again, is often implied by the problem: when $G =$ Floridians, $-G =$ Rest of US population; when $G =$ African Americans, $-G =$ White Americans. A distinctive prediction of our model, confirmed by our experiments in Section 4.1, is that the stereotype for a given group G depends on the comparison group $-G$.¹¹ When $-G$ is not pinned down by the problem itself, to derive testable predictions from representativeness, we set $-G = \Omega \setminus G$ where Ω is the natural population over which the unconditional distribution of types is measured.

3 Properties of Stereotypes

We now study stereotypical beliefs and their accuracy. To illustrate the role of representativeness, we first ask to what extent the most representative type is a good fit for the group, namely whether it is modal. Next, we assess the accuracy of the entire stereotypical distribution. To do so, we focus on a cardinal types and compare the stereotype’s mean and

¹⁰For example, suppose a person is asked to guess the typical occupation of a democratic voter in an “open ended” format (without being provided with a set of alternatives). Here the level of granularity at which types are defined is not obvious (e.g. teacher vs a university teacher vs a professor of comparative literature).

¹¹Some empirical papers have taken a similar approach, exogenously varying the natural comparison group through priming. Benjamin, Choi, and Strickland (2009) show that priming racial or ethnic identity can impact the risk preferences of participants. Chen et al (2010) find that Asian students cooperate less with outgroup members when primed with their ethnic identity rather than their university identity. Shih et al (1999) show that Asian-American women self-stereotype themselves as better or worse in math, with corresponding impact on performance, when their ethnicity or gender, respectively, is primed. Shih et al (2006) replicate this effect using a verbal task, documenting that Asian-American women performed better when their gender rather than their ethnicity was primed. While the generalizability and replicability of priming has been doubted (Klein et al 2014), this body of evidence is consistent with context dependence.

variance to the true ones.

3.1 Likely vs Unlikely Exemplars

The most representative type for a group is the one that agents most easily recall and associate with the group itself, for instance a red-haired Irishman or a 65+ year old Floridian. Social psychologists call this type the exemplar of the group. Under any specification of the weighting function h_t in Definition 2, overweighting (weakly) increases as we move toward more representative types, so the exemplar is also the type whose probability is overweighted the most.¹² By analyzing the exemplar, then, we can gauge whether representativeness induces the DM to overweight a likely type (as it happens standard models of categorical thinking) or an unlikely type. When overweighting occurs in unlikely and extreme types, the biases of stereotypes can be particularly severe.¹³

Equation (2) then yields the following characterization.

Proposition 1 *Suppose the conditional distributions π_G and π_{-G} are not identical. Consider two extreme cases:*

i) If for all $t, t' \in T$ we have that $\pi_{t,G} > \pi_{t',G}$ if and only if $\pi_{t,-G} > \pi_{t',-G}$, then the exemplar is not the modal type for at least one group.

ii) If for all $t, t' \in T$ we have that $\pi_{t,G} > \pi_{t',G}$ if and only if $\pi_{t,-G} < \pi_{t',-G}$, then for each group the exemplar is the modal type.

Case i) says that when groups have similar distributions, in the sense of having the same likelihood ranking, the most representative type is unlikely for at least one group, potentially for both. Representativeness draws the DM's attention to group differences, neglecting the fact that the groups are similar, and have the same mode. This mechanism generates inaccurate stereotypes and is illustrated by the Florida example.

¹²Consider the function $h(\cdot)$ from Definition 2. When applied to more more representative types, the first argument of the function increases, while one of the other $T - 1$ arguments decreases. As a result the weighting factor h_t (and thus overweighting $\pi_{t,G}^{st}/\pi_{t,G}$) increases as well.

¹³In the rank-based truncation model, the frequency of the exemplar provides a measure of stereotype accuracy. By accuracy, we mean the extent to which the stereotype minimizes the distance $\sum_t (\pi_{t,G}^{st} - \pi_{t,G})^2$. When $d = 1$ and only one type is recalled (there are no ties), accuracy is maximized if the exemplar is the most likely type and minimized if the exemplar is the least likely type.

Case ii) says that the most representative type tends to be likely for both groups when the distributions are very different. In this case, groups differ the most around their modes, so representativeness and likelihood coincide. Thinking of Swedes as “blond haired” and Europeans as “dark haired” is accurate precisely because these are majority traits of the Swedish and European populations, respectively. In these cases, stereotyping yields fairly reliable models. Of course, there is still some inaccuracy. Even in the case of likely exemplars, judgment errors can be significant. For instance, voters in some U.S. states are perceived as “blue” or “red” because a majority of the population indeed votes Democrat or Republican. In reality, even in “blue” states, far from everyone votes Democrat. In the 2012 Presidential election, vote shares of either candidate in most states ranged from 40% to 60%.¹⁴

When DMs strongly overweight representative types, the most severe biases occur when those types are unlikely and extreme. This is true both under rank based truncations and under smooth discounting functions (see Section 3.2). Ethnic stereotypes based on crime or terrorism exhibit this error: they neglect the fact that by far the most common types in all groups are honest and peaceful.

3.2 Stereotypical Moments

We now characterize how the first two moments of a distribution are distorted by the process of stereotyping. To do so, we must restrict our analysis to cardinal, ordered types. The following results hold for any weighting function $h_t(\cdot)$ satisfying Definition 2. We consider two canonical cases that prove useful in illustrating the predictions of the model.

In the first case, groups G and $-G$ are such that the likelihood ratio $\pi_{t,G}/\pi_{t,-G}$ is monotonic in t . The monotone likelihood ratio property (MLRP) holds to a first approximation in many empirical settings and is also assumed in many economic models, such as standard agency models.¹⁵ If $\pi_{t,G}/\pi_{t,-G}$ is monotonically increasing (decreasing) in t , then group G is

¹⁴See https://en.wikipedia.org/wiki/United_States_presidential_election,_2012, section on votes by electoral college.

¹⁵Examples include the Binomial and the Poisson families of distributions with different parameters. The characterisation of distributions satisfying MLRP is easier in the case of continuous distributions, see Appendix D: two distributions $f(x)$, $f(x - \theta)$ that differ only in their mean satisfy MLRP if and only if the distribution $f(x)$ is log-concave. Examples include the Exponential and Normal distributions. To the extent that discrete distributions sufficiently approximate these distributions (as the Poisson distribution $Pois(\lambda)$ approximates the Normal distribution $N(\lambda, \lambda)$ for large λ), they will also satisfy MLRP.

associated with higher (lower) values of t relative to the comparison group $-G$. Formally:

Proposition 2 *Suppose that MLRP holds. Then, for any weighting function $h_t(\cdot)$:*

i) If the likelihood ratio $\frac{\pi_{t,G}}{\pi_{t,-G}}$ is strictly increasing in t , then

$$\mathbb{E}^{st}(t|G) > \mathbb{E}(t|G) > \mathbb{E}(t|-G) > \mathbb{E}^{st}(t|-G),$$

ii) If the likelihood ratio $\frac{\pi_{t,G}}{\pi_{t,-G}}$ is strictly decreasing in t , then

$$\mathbb{E}^{st}(t|G) < \mathbb{E}(t|G) < \mathbb{E}(t|-G) < \mathbb{E}^{st}(t|-G).$$

Under MLRP, the most representative part of the distribution for G is the right tail if $\pi_{t,G}/\pi_{t,-G}$ increases in t or the left tail if $\pi_{t,G}/\pi_{t,-G}$ decreases in t . The representative tail is then overweighted while the non-representative tail is underweighted. As a consequence, the assessed mean $\mathbb{E}^{st}(t|G)$ is too extreme in the direction of the representative tail.

Critically, in line with the social cognition perspective, the stereotype contains a kernel of truth: the DM overestimates the mean of G if this group has a higher mean than the comparison group, namely $\mathbb{E}(t|G) > \mathbb{E}(t|-G)$ and conversely if $\mathbb{E}(t|G) < \mathbb{E}(t|-G)$. The DM exaggerates this true difference because he inflates the association of G with its most representative types.¹⁶ For instance, when judging an asset manager who performs well, we tend to over-emphasize skill relative to luck because higher skill levels are relatively more associated with higher performance. This occurs even if for both skilled and unskilled managers high performance is mostly due to luck.

In the second case for which we characterize the stereotypical distributions, groups G and $-G$ have the same mean $\mathbb{E}(t|G) = \mathbb{E}(t|-G) = \mathbb{E}(t)$ but differ in their variance. We abstract from skewness and higher moments by considering distributions $(\pi_{t,G})_{t \in T}$ and $(\pi_{t,-G})_{t \in T}$ that share the same support and are both symmetric around the median/mean $\mathbb{E}(t)$.

¹⁶Depending on the distribution and the weighing function, the DM's assessment of the variance $\text{Var}(t|G)$ may also be dampened relative to the truth. This is often true under the truncation weighing function. In this case, stereotyping effectively leads to a form of overconfidence in which the DM both holds extreme views and overestimates the precision of his assessment. That extreme views and overconfidence (in the sense of over precision) go together has been documented in the setting of political ideology, among others (Ortoleva and Snowberg 2015).

Proposition 3 *Suppose that in G more extreme types are relatively more frequent than in $-G$. Formally, the likelihood ratio $\frac{\pi_{t,G}}{\pi_{t,-G}}$ is U-shaped in t around $\mathbb{E}(t)$. Then, for any weighting function $h_t(\cdot)$ stereotypical beliefs satisfy:*

$$\begin{aligned} \text{Var}^{st}(t|G) &> \text{Var}(t|G) > \text{Var}(t|-G) > \text{Var}^{st}(t|-G), \\ \mathbb{E}^{st}(t|G) &= \mathbb{E}^{st}(t|-G) = \mathbb{E}(t). \end{aligned}$$

When group G has a higher relative prevalence of extreme types, its representative types are located at both extremes of the distribution. The DM’s beliefs about G are then formed by overweighting both tails while underweighting the unrepresentative middle. The overweighting of G ’s tails causes the assessment of its variance $\text{Var}^{st}(t|G)$ to be too high. For example, the skill distribution of immigrants to the US may be perceived as having very fat tails, or even bimodal, with immigrants being perceived as either unskilled or very skilled relative to the native population. The mean of the group, in contrast, is assessed correctly, because the stereotypical distribution remains symmetric around $\mathbb{E}(t)$. As before, the stereotype contains a kernel of truth. It induces the agent to exaggerate the true differences between groups, namely the higher variance of G relative to its counterpart.

We present a number of extensions of the model in the Online Appendix. We first consider multi-dimensional type spaces, and show that stereotypes center around the dimension where groups differ the most, in line with the kernel of truth logic (Online Appendix C). Multidimensional stereotypes imply that the dimension we think about is influenced by context dependent. For example, the Irish are stereotyped as red-haired when compared to the European population. However, when compared to the Scots, a more plausible stereotype for the Irish is “Catholic” because religion is the dimension along which Irish and Scots differ the most.

In Online Appendix D we extend the model to continuous type spaces. Many settings of interest in economics can be usefully described by continuous probability distributions, and we show our model is particularly tractable in this case. In Online Appendix E, we relax Definition 2 and allow weighting of types to also be influenced by their likelihood. We show that the basic insight that stereotypes contain a kernel of truth carries through to each of

these cases as well.

Finally, in Online Appendix F we embed our model into a learning setup in which information arrives over time. In this setting we explore how stereotypical thinking distorts reaction to information. So long as stereotypes do not change, people under-react or even ignore information inconsistent with stereotypes. If enough contrary information is received, stereotypes change, leading to a drastic reevaluation of already available data. Representativeness-based recall reconciles under-reaction with over-reaction to data, generating both confirmation bias and base-rate neglect (Gennaioli, Shleifer, and Vishny 2015).

To summarise, the psychology of representativeness yields stereotypes that are consistent with the social cognition approach in which individuals assess groups by recalling and focusing on distinctive group traits. When there are systematic differences between groups, stereotypes get the direction right, but exaggerate differences.

3.3 Some Examples

A growing body of field and experimental evidence points to a widespread belief that women are worse than men at mathematics (Eccles, Jacobs, and Harold 1990, Guiso, Monte, Sapienza and Zingales 2008, Carrell, Page and West 2010). This belief persists despite the fact that, for decades, women have been gaining ground in average school grades, including mathematics, and have recently surpassed men in overall school performance (Goldin, Katz and Kuziemko 2006, Hyde et al, 2008). This belief, shared by both men and women (Reuben, Sapienza and Zingales 2014), may help account, in part, for the gender gap in the choices of high school tracks, of college degrees and of careers, with women disproportionately choosing humanities and health related areas (Weinberger 2005, Buser, Niederle and Oosterbeek 2014) and foregoing significant wage premiums to quantitative skills (Bertrand 2011).

Gender stereotypes in mathematics, particularly beliefs that exaggerate the extent of average differences, are consistent with the predictions of our model. The fact that men are over-represented at the very highest performance levels leads a stereotypical thinker to exaggerate the magnitude of mean differences. Figure 1 shows the score distributions from the mathematics section of 2013's Scholastic Aptitude Test (SAT), for both men and

women.¹⁷ The distributions are very similar, with average scores being slightly higher for men (531 versus 499 out of 800). However, scores for men have a heavier right tail, with men twice as likely to have a perfect SAT math score than women.¹⁸ In light of such data, the stereotypical male performance in mathematics is high, while the stereotypical female performance is poor. Predictions based on such stereotypes are inaccurate, exaggerating true differences. While differences in the right tail of the distribution are unlikely to be relevant for most decisions, stereotypical thinking driven by these differences has the potential to impact economically-important decisions, whether through self-stereotyping (i.e., choice of careers or majors as in Buser, Niederle, Osterbeek 2014) or through discrimination (i.e., hiring decisions as in Bohnet, van Geen, and Bazerman 2015).

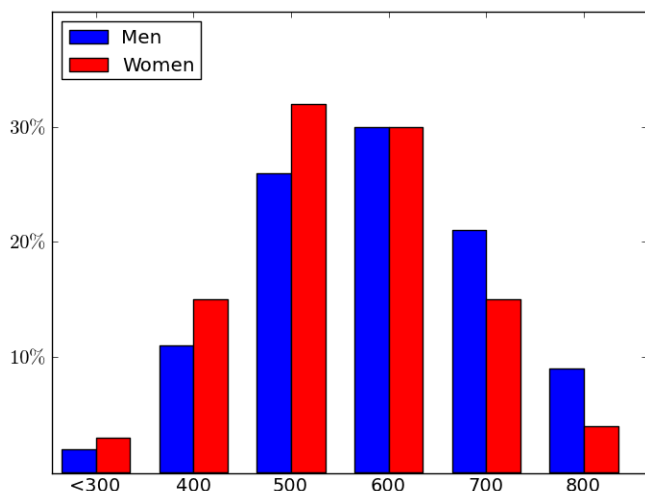


Figure 1: SAT Mathematics scores by gender (2012)

The logic of exaggerated, yet directionally correct, stereotypes can also shed light on the well documented phenomenon of base rate neglect (Kahneman and Tversky, 1973). Indeed,

¹⁷Standardized test performance measures not only innate ability but also effort and investment by third parties, Hyde et al (2008). The mapping of test performance into inferences about innate ability is an issue not addressed by our model.

¹⁸For 2013 SAT Mathematics scores, see <http://media.collegeboard.com/digitalServices/pdf/research/SAT-Percentile-Ranks-By-Gender-Ethnicity-2013.pdf>. Results are similar for the National Assessment of Educational Progress (NAEP), which are more representative of the overall population. For 2012 NAEP scores for 17 year olds in mathematics, see http://nationsreportcard.gov/ltt_2012/age17m.aspx. See Hyde et al (2008), Fryer and Levitt (2010), and Pope and Sydnor (2010) for in-depth empirical analyses of the gender gap in mathematics.

Proposition 2 implies that the DM overreacts to information that assigns people to groups, precisely because such information generates extreme stereotypes.¹⁹ Consider the classic example in which a medical test for a particular disease with a 5% prevalence has a 90% rate of true positives and a 5% rate of false positives. The test assigns each person to one of two groups, + (positive test) or - (negative test). The DM estimates the frequency of the sick type (s) and the healthy type (h) in each group. The test is informative: a positive result increases the relative likelihood of sickness, and a negative result increases the relative likelihood of health for any prior. Formally:

$$\frac{\Pr(+|s)}{\Pr(+|h)} > 1 > \frac{\Pr(-|s)}{\Pr(-|h)}. \quad (4)$$

This condition has clear implications: the representative person who tests positive is sick, while the representative person who tests negative is healthy. Following Proposition 2, the DM reacts to the test by moving his priors too far in the right direction, generating extreme stereotypes. He greatly boosts his assessment that a positively tested person is sick, but also that a negatively tested person is healthy. Because most people are healthy, the DM's assessment about the group that tested negative is fairly accurate but is severely biased for the group that tested positive. This analysis formalises Tversky and Kahneman's (1983) verbal account of base rate neglect.²⁰

Exaggerated stereotypes may shed light on several other phenomena. When assessing the performance of firms in a hot sector of the economy, the investor recalls highly successful (and some moderately successful) firms in that sector. However, he neglects the possibility of failures, because failure is statistically non-diagnostic, and psychologically non-

¹⁹In Online Appendix F we explore in detail how stereotypical beliefs react to a different kind of information, namely information about the distribution of types when groups are *given*.

²⁰Our account is distinct from a mechanical underweighting of base-rates in Bayes rule, as in Grether (1980) and Bodoh-Creed, Benjamin and Rabin (2013). In those models, upon receiving the test results, the DM can update his beliefs in the wrong direction: he can be less confident that a person is healthy after a negative test than under his prior, which cannot happen in our model.

While this prediction of our model seems consistent with introspection, we are not aware of experimental evidence on this point. Griffin and Tversky (1992) present evidence consistent with pure neglect of base rates, but in a significantly different task, namely inferring the bias of a coin from a history of coin flips. Such experiments are hard to compare with the predictions of our model, because subjects are asked to generate distributions of different numbers of coin flips in their minds, which is a much more involved task than to recall types of a given distribution. Their assessments, then, might be wrong for other reasons. See Bodoh-Creed, Benjamin and Rabin (2013) for a detailed discussion.

representative, of a growing sector – even if it is likely. This causes both excessive optimism (in that the expectation of growth is unreasonably high) and overconfidence (in that the variability in earnings growth considered possible is truncated). True, the hot sector may have better growth opportunities on average, but representativeness exaggerates this feature and induces the investor to neglect a significant risk of failure. Similarly, when assessing an employee’s skill level, an employer attributes high performance to high skill, because high performance is the distinctive mark of a talented employee. Because he neglects the possibility that some talented employees perform poorly and that some non-talented ones perform well (perhaps due to stochasticity in the environment), the employer has too much faith in skill, and neglects the role of luck in accounting for the output.

4 Evidence on Representativeness and Stereotypes

In testing our model, we focus on the two main implications of representativeness-based stereotypes:

- **Kernel of truth:** stereotypes depend on group characteristics, and – in most (precisely characterized) settings – are slanted toward representative types.
- **Context dependence:** the stereotype of a target group depends on the characteristics of the reference group it is compared to.

We first test these properties with a lab experiment (Section 4.1). We then turn to survey data on beliefs about U.S. political groups (Section 4.2) for an empirical analysis. The survey data is more tightly linked to our interest in social stereotypes. The laboratory experiment, however, allows us to directly test the role of representativeness in generating context dependent beliefs. Online Appendices G and H provide all details, and additional results, for the experiments and field evidence.

4.1 Lab Evidence on Representativeness and Context Dependence

The influence of the representativeness heuristic on recall and on beliefs has been extensively documented in the lab (Kahneman and Tversky 1972, 1983). Our goal here is to consider

how representativeness as formalized in Equation (2) gives rise to context dependent beliefs. To our knowledge, the possibility that representativeness may generate context dependence has not been tested before.

To assess this prediction, we perform a controlled laboratory experiment that allows us to isolate representativeness from many confounding factors – historical, sociological, or otherwise – that may affect stereotype formation in the real world. We construct our own groups of ordinary objects, creating a target group, G , and a comparison group, $-G$. We hold the target group G fixed, but explore how participant impressions of it change as we change the comparison group $-G$, and hence representativeness.

We conducted several experiments, in the laboratory as well as on Amazon Mechanical Turk. Each involves a basic three-step design. First, participants are shown the target group and a randomly-assigned comparison group for 15 seconds. In this time, differences between groups can be noticed but the groups’ precise compositions cannot be memorized. The second step consists of a few filler questions, which briefly draw the participants’ cognitive bandwidth away from their observation. Finally, participants are asked to assess the groups they saw. Participants are incentivized to provide accurate answers.

We randomly assign participants to either the Control or the Representativeness condition. In the Control condition, G and $-G$ have nearly identical distributions, so that all types are equally representative for each group. In the Representativeness condition, $-G$ is changed in such a way that a certain type becomes very representative for G . Context dependence implies that the assessment of G should now overweight this representative type, even though the distribution of G itself has not changed.

We ran six experiments of this form, with design changes focused on reducing participant confusion and removing confounds. Here, we describe the final, and most refined, version of these experiments. In an attempt to provide a overview of the results while remaining concise, we also provide the results from pooled specifications that use all data collected. In Appendix G, we present additional details and report all experiments conducted. We also provide instructions and materials for each experiment and the full data set.

Consider first the experiment illustrated in Figure 2. A group of 25 cartoon girls is presented next to a group of 25 cartoon boys in t-shirts of different colors: blue, green,

or purple. In the Control condition, Fig.2a, the groups have identical color distributions (13 purple, 12 green), so no color is representative of either group. The Representativeness condition, Fig.2b, compares the *same* group of girls with a different group of boys, for whom green shirts are replaced by blue shirts. Now only girls wear green and only boys wear blue. These colors, while still not the most frequent for either group, are now most representative. For each group, girls and boys, participants are asked a number of questions concerning the frequency of T-shirts of different colors worn by that group.

Applying our model, in the control condition the type space is $T = \{\text{green, purple}\}$, and the groups are $G = \text{girls}$, and $-G = \text{boys}$. Given that the color distributions are identical across groups, both types are equally representative, $\pi_{\text{green,girls}}/\pi_{\text{green,boys}} = 1 = \pi_{\text{purple,girls}}/\pi_{\text{purple,boys}}$. As a result, assessment of G should be on average correct, $\pi_{\text{green,girls}}^{\text{st,control}} = \pi_{\text{green,girls}}$ and $\pi_{\text{purple,girls}}^{\text{st,control}} = \pi_{\text{purple,girls}}$ for any weighing function (and the same is true about assessments of $-G$).

In the treatment condition, the distribution of shirt colors remains the same for girls. For boys, green shirts are changed into blue. Thus, the type space changes to $T = \{\text{green, purple, blue}\}$ and the representative color for girls becomes green, $\pi_{\text{green,girls}}/\pi_{\text{green,boys}} = \infty > 1 = \pi_{\text{purple,girls}}/\pi_{\text{purple,boys}}$, while that for boys becomes blue. As a result, in the treatment condition subjects should inflate the frequency of green shirts relative to the truth, $\pi_{\text{green,girls}}^{\text{st,treatment}} > \pi_{\text{green,girls}}$ (and the same should happen to assessments of blue shirts for boys). We also expect the assessed frequency of green shirts to go up relative to the control condition, namely $\pi_{\text{green,girls}}^{\text{st,treatment}} > \pi_{\text{green,girls}}^{\text{st,control}}$. Critically, the only factor that varies across treatments is the representativeness of the 12-color shirt. Thus, if we see differences across conditions, the causal role of representativeness-based recall in shaping group judgments is clear.²¹

We collected data from 301 participants using this T-shirts design.²² Since the number of green and purple shirts is very similar, we first ask subject the simplest question of which shirt color is modal. Next, we ask subjects to assess the share of green and purple shirts.

²¹We vary which colors are used in which roles across participants. Some participants saw this particular color distribution, while others see, for example, green as the modal color, with purple as the diagnostic color for boys in the Rep. condition and blue as the diagnostic color for girls in the Rep. condition. We vary the colors across the roles to avoid confounding the characteristics of any particular color with its diagnosticity.

²²Throughout our analysis, we exclude any participant who participated in a previous version of the experiment and any participant who self-identified as color blind. In Appendix G, we show that our results are unchanged if we include these additional observations.

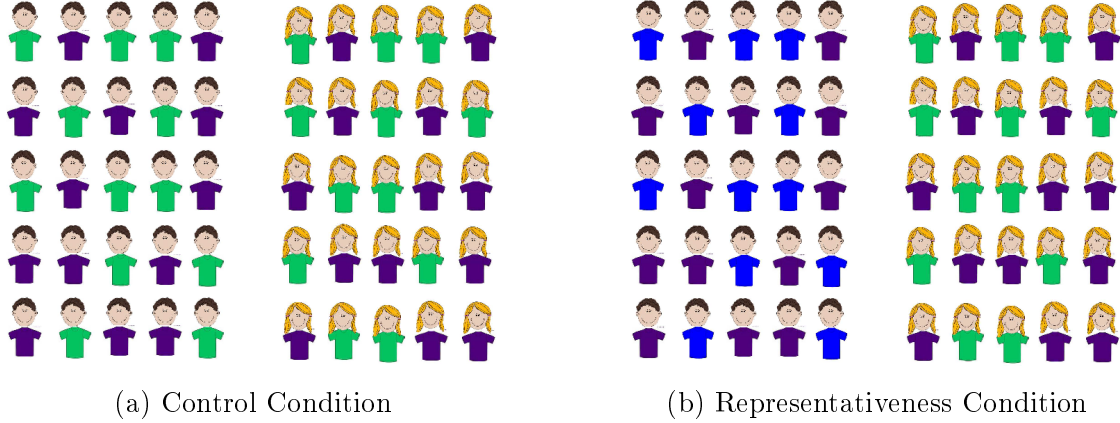


Figure 2: T-shirts Experiment

Consistent with the role of representativeness, participants assigned to the Representativeness condition are 10.5 percentage points more likely to recall the less frequent color, green for girls or blue for boys, as the modal color when it is representative of a group (35% of participants guess the less frequent color is modal in the Control condition, this proportion increases to 46% in the Representativeness condition, $p=0.01$, estimated from a probit regression reported in Appendix G).

Let us now turn to subjects' estimates of how many T-shirts of each color they saw in each group. In both conditions, the true difference in counts is one (13 purple shirts, 12 green or blue shirts). In the Control condition, participants on average believe they saw 0.54 more purple shirts than green or blue shirts. In the Representativeness condition, participants believe they saw 0.72 fewer purple shirts than green or blue shirts (the across treatment difference is significant with $p=0.013$ from two-tailed Fisher Pitman permutation test).

In total, we collected data for six experiments of this general structure, gathering evidence from more than 1,000 participants. As we describe in Appendix G, while there is substantial variation across experiments, when we pool all data collected we find significant aggregate treatment effects in line with a role of representativeness in judgment.²³ Using a probit regression that pools all of the data for unordered type experiments similar to the T-shirts experiment (four versions, 741 participants), we find that participants are 9.3 percentage points more likely to guess that the less frequent type is modal when it is representative than

²³We find effects in the predicted directions for all six designs, with significant effects for two. We discuss the extent to which our results are sensitive to the specifics of the design in Appendix G.

when it is not ($p=0.002$). We also run a family of ordered types experiments (two versions, 402 participants). Overall, in those experiments, participants are 9.3 percentage points more likely to guess that the group of interest has a greater average than the comparison group when the right tail is representative ($p=0.062$).²⁴ Given our simple experimental setting with groups of mundane objects, we interpret our results – a significant and reasonably-sized impact on average beliefs – as an important proof of concept: the presence of representative types biases *ex post* assessment.

4.2 Empirical Evidence on Political Stereotypes

We examine two data sets on political preferences, and beliefs about political preferences, in the U.S. We investigate the roles of representativeness and context dependence by separately testing for hypotheses that allow us to assess the leading theories of stereotypes.

First, we test whether beliefs are correct or depart systematically from the truth. The statistical discrimination approach builds on the assumption that people hold rational expectations of group traits. Comparing beliefs to the truth allows us to assess the validity of this assumption in our data.

Second, we test if beliefs depart from the truth by exaggerating (mean) differences among groups, as per the kernel of truth hypothesis. This is an implicit test of context dependence, because it implies that beliefs about the target groups are shaped not only by that group’s characteristics, but also by those of the reference group.²⁵

Third, we test if distortions in beliefs can be accounted for by the overweighting of highly representative types (defined as types that are relatively more frequent in the target relative to the reference group). The second and third tests address the key predictions of our model.

Finally, models of categorization (Mullainthan 2003, Fryer and Jackson 2008) predict that beliefs about average group traits can be distorted if they exaggerate the incidence of the most likely group type. At the end of this section we present a test to discriminate among

²⁴Results for the ordered types experiments, unlike the simpler T-shirts style design, were sensitive to the choice of platform, with consistently strong results on Amazon Mechanical Turk and weak or null results in laboratory samples. We discuss this in Appendix G.

²⁵Of course, unlike in the laboratory experiment, in this setting we cannot test for context dependence by exogenously varying the reference group.

likelihood-based stereotypes and the representativeness-based stereotypes that we propose.

4.2.1 The data

We have two data sets on political preferences and beliefs about political preferences. The first data set, from Graham et al (GNH 2012), contains data from the Moral Foundations Questionnaire. Respondents (1,174 self-identified liberals and 500 self-identified conservatives) answer questions about their position on a subset of 45 issues: 20 moral relevance statements (e.g., “when you decide whether something is right or wrong, to what extent are the following considerations relevant to your thinking?”) and 25 moral judgments (e.g., “indicate the extent to which you would agree or disagree”). For each issue, a randomly determined subset of participants states their own position, another subset states their belief on the position of a “typical liberal”, and a third subset states their belief on the position of a “typical conservative”. The data thus includes the distribution over positions for both liberal and conservatives, as well as the average believed typical position of liberals and of conservatives, on each of the 45 issues. Each position is elicited on 1 - 6 scale.

The second data set comes from the American National Election Survey (ANES), and contains data from more than 20,000 respondents between 1964 and 2012.²⁶ The survey covers political issues of the day, such as the optimal amount of government spending and service provision (1984 through 2000), or the proper place of women in society (1972 through 1998). We focus on the 10 issues that ask participants to respond on a multi-point, 1 to 7, scale (rather than just indicate binary agreement or disagreement); each of these 10 issues is asked in multiple years. Participants are asked to provide their own position on the scale and their believed position of the Democratic and Republican party (“Where would you place the Democratic (Republican) party on this scale?”). The data includes, for each issue-year observation, the distribution of participant positions for both self-identified liberals and self-identified conservatives, as well as the distribution of believed typical positions of the Democratic and Republican Parties.

²⁶This data is publicly available at http://www.electionstudies.org/studypages/anes_timeseries_cdf/anes_timeseries_cdf.htm.

4.2.2 Empirical strategy and results

Our analysis focuses on beliefs about two groups, Conservatives and Liberals. The types are the possible positions for each issue (1, 2, ..., 6, 7). For the GNH data, we interpret beliefs about the “typical” element of a group to coincide with the believed average position in that group. Similarly, for the ANES data we use the believed party positions as a proxy for believed mean of each group.²⁷ We then take as a benchmark the hypothesis that individuals hold accurate beliefs about each group, and in particular that believed mean position should equal true mean position, at least on average across subjects. The accurate beliefs hypothesis underlies the most common economic model of stereotyping, statistical discrimination.

To assess our representativeness-based model, we perform a regression exercise. Under the assumption that the weighting function $h_t(\cdot)$ is differentiable, our model yields two regression specifications that we estimate.

Proposition 4 *Let $G \in \{\text{conservative, liberal}\}$, and let $h_t \equiv h(\pi_{t,G}/\pi_{t,-G})$ be a differentiable weighting function as in Definition 2. We then have:*

1) *Kernel of truth regression. There exists a constant $\kappa > 0$ such that, as a first order approximation around identical distributions $\pi_G/\pi_{-G} = \mathbf{1}$, we have:*

$$\mathbb{E}^{st}(t|G) = \mathbb{E}(t|G) (1 + \kappa) - \kappa \cdot \mathbb{E}(t| - G). \quad (5)$$

2) *Representativeness regression. Denote $H = \{T - 2, \dots, T\}$ the right tail of types and $\Theta = \sum_H \pi_{t,cons} / \sum_H \pi_{t,lib}$ as the average representativeness of right tail types for conservatives. Our model entails the following approximate equations:*

$$\mathbb{E}^{st}(t|cons) = \mathbb{E}(t|cons) + \lambda_{cons} (\Theta - 1), \quad (6)$$

$$\mathbb{E}^{st}(t|lib) = \mathbb{E}(t|lib) - \lambda_{lib} (\Theta - 1), \quad (7)$$

²⁷This assumption is consistent with the authors’ interpretation of the GNH data (GNH 2012) and with previous studies using ANES (e.g., Westfall et al, WBCJ 2015). Furthermore, to the extent that this assumption holds equally well for most issues within a data set, our focus on across-issue differences should allow us to test the predictions of our model even with an imperfect proxy for beliefs of mean positions. Finally, the data provides some insight into whether subjects are reporting (perceived) modal or mean types. As we show below, the modal type is a poor prediction of stated beliefs, while a distorted mean slanted towards representative types is an accurate prediction of stated beliefs.

where λ_{cons} and λ_{lib} are positive constants that depend on the true distributions π_{cons} and π_{lib} .

The first regression allows us to test for the kernel of truth hypothesis, while the second set of regressions allows us to test for the role of representativeness.

Equation (5) says that respondents in our model inflate the average position of a group, say the conservatives, if and only if the group has a higher average position than the other group, namely the liberals. Formally, $\mathbb{E}^{st}(t|cons) > \mathbb{E}(t|cons)$ if and only if $\mathbb{E}(t|cons) > \mathbb{E}(t|lib)$. Because in our measurement scale higher types mean “more conservative”, we expect: i) believed conservative average to be higher than the truth, and ii) the extent of overstatement to decrease in the average liberal position $\mathbb{E}(t|lib)$. Conversely, we expect the average liberal position to be lower than the truth, the more so the higher the average conservative position $\mathbb{E}(t|cons)$.

As previously discussed, the basis of these predictions is context dependence: information about the distribution of $-G$ is relevant for the beliefs about G . This context dependence is inconsistent with rational expectations, in which only the group’s own means should affect beliefs. We test the hypothesis that the true mean $\mathbb{E}(t|G)$ is a significant predictor of the believed mean $\mathbb{E}^{st}(t|G)$ with a positive sign, while the other group’s true mean $\mathbb{E}(t|-G)$ is a predictor of the believed mean with a negative sign.

Equations (6) and (7) say that respondents’ assessment bias is shaped by representativeness. When the right tail is more representative for conservatives, $\Theta > 1$, participants should inflate the average conservative position more (higher $\mathbb{E}^{st}(t|cons) - \mathbb{E}(t|cons)$) and deflate the average liberal position more (lower $\mathbb{E}^{st}(t|lib) - \mathbb{E}(t|lib)$). We test the hypothesis that the inflation in conservative positions is positively associated with the representativeness of the right tail for the conservatives, while the inflation in liberal positions is negatively associated with it. Once again, the representativeness of the right tail is computed using the true distribution of positions.

In Equations (6) and (7), in many cases the representative tail is also the most likely one. As a consequence, these tests cannot distinguish a representativeness-based from a likelihood-based model of distorted beliefs. We perform two additional tests. First, we run versions of Equations (6) and (7) in which we control for the likelihood of tails (see Table 6 in

Appendix H). Second, we compute numerically the predictions of a representativeness-based model of stereotypes and of a likelihood-based model of stereotypes. We then assess which of these two is better able to match the data on beliefs.

4.2.3 Empirical Results

To begin, we illustrate the structure of the data and the nature of our predictions with two simple examples from the GNH data set, focusing on beliefs about conservatives. In Example 1, participants are asked about their agreement with the statement, “It can never be right to kill a human being”. In Example 2, participants are asked about the moral relevance of “whether or not someone cared for someone weak or vulnerable”. As can be seen in Figure 3, in Example 1 the modal position (Strongly Disagree (1)) and most representative positions (Strongly Disagree (1)) coincide for conservatives. In contrast, in Example 2 in Figure 3, the most representative types (Slightly Relevant (3), Not at all Relevant (1)) are not most likely for the conservative group. Following Proposition 1, we predict that beliefs will be distorted in the direction of the most representative types. Thus, we expect more exaggeration in Example 2 than in Example 1, since in Example 2 the most representative types (in the left tail) are far from the modal type, while in Example 1, they coincide. This is what we find: the conservative position is exaggerated by only 0.09 positions in Example 1 (true mean 2.99, believed mean 2.90), but by 1.06 positions in Example 2 (true mean 4.21, believed mean 3.15).

In the full data sets, we treat each (issue, year) pair as an observation, and we cluster standard errors at the issue level. For the GNH data, we have 45 observations: 45 issues each measured in the same year. For the ANES data, we have 66 observations: 10 issues, each measured in multiple years. To begin, we simply document systematic exaggeration in both data sets. This is a primary focus of the original analysis in GNH (2012), and also in WBCJ (2015)’s analysis of the ANES data. Figure 4 shows that the believed difference between typical conservative and typical liberal positions is larger than the true difference in mean positions for 109 of the 111 observations. The data for both GNH (purple squares) and ANES (orange triangles) lie above the 45 degree line (dashed).²⁸ Average exaggeration

²⁸For convenience, we recode all issues so that the high end of the scale (6,7) represents the stereotypically

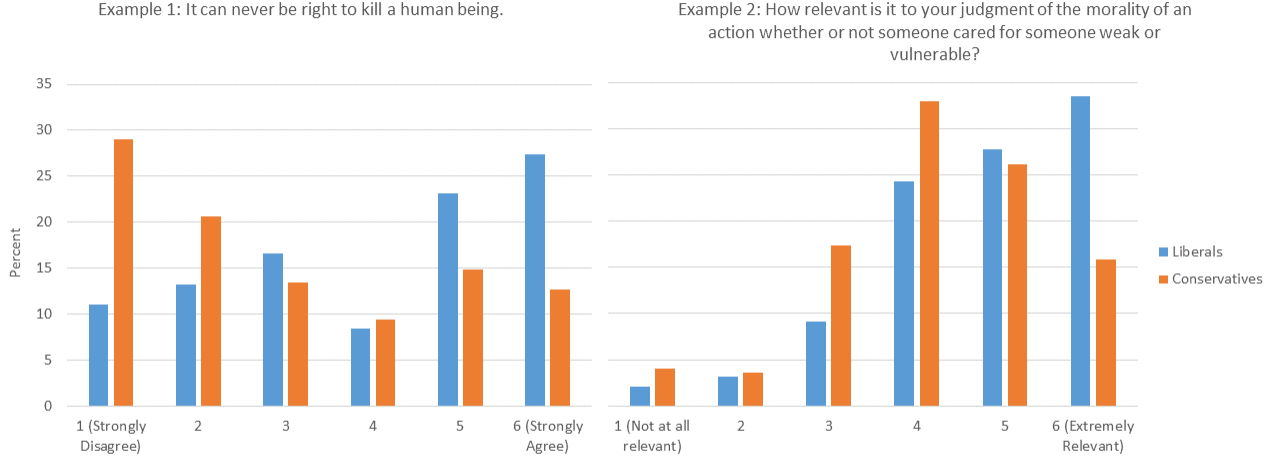


Figure 3: Two Examples

is 0.62 positions on the scale (0.66 in the GNH data, 0.59 in the ANES data).²⁹

The systematic and significant exaggeration of mean differences suggests that the benchmark model of accurate beliefs is missing something important. Indeed, this exaggeration reflects the fact that believed means are typically more extreme than true means. First, note that the Kernel of truth regression (Equation 5) generates exaggeration of mean differences, just as documented in Figure 4:

$$\mathbf{E}^{st}(\mathbf{t}|\mathbf{G}) - \mathbf{E}^{st}(\mathbf{t}|- \mathbf{G}) = (\mathbf{1} + \mathbf{2}\kappa) \cdot [\mathbb{E}(t|G) - \mathbb{E}(t|-G)]$$

The results from regression (5) are shown in Table 1. While these results provide strong evidence of context dependence and are consistent with our model, they do not pin down a role for representativeness of types. Our next test relates the magnitude of representativeness

more conservative position.

²⁹A natural question to ask is how beliefs vary across liberals and conservatives. That is, do beliefs about a group G depend on membership in G versus $-G$. Our model does not speak to this issue. However, in Appendix H, we show that the results we document below hold for both beliefs held by conservatives and beliefs held by liberals; see Tables 11, 12, and 13.



Figure 4: Exaggeration of Differences

of tail types to the magnitude of belief distortions.

To this end, we implement the regressions in Equations (6, 7). Following Proposition 4, we compute the average representativeness of tail types for conservatives, $\Theta = \frac{\sum_{t>T-2} \pi_{t,cons}}{\sum_{t \geq T-2} \pi_{t,lib}}$. We again test the hypothesis that Θ is a significant predictor of $\mathbb{E}^{st}(t|cons)$ with a positive sign, and a predictor of $\mathbb{E}^{st}(t|lib)$ with a negative sign. Table 2 shows that, conditional on true mean, Θ predicts believed mean for each group G as predicted. In Appendix H Table 6, we repeat this analysis and obtain similar results controlling for the average likelihood on the tail positions. This rules out that what drives these effects is that representative tails are also likely tails, and indicates the additional effect of tails being representative.

Finally, we use the model to predict the believed mean of group G . For simplicity, we focus on the rank-based truncation specification, which is easiest to apply to the data. We then assess predictions about mean beliefs when stereotypes include only the d most representative types, for $d = 1, \dots, T$. We compare these predictions to those of a model in which beliefs are obtained by restricting the distribution to the d most likely types. Our benchmark for both models is predicting the believed mean from the entire distribution, where $d = T$.

Table 1: Information about -G Predicts Beliefs about G

| OLS Predicting Believed Mean of Group G | | | | | | |
|---|---------------------|----------------------|----------------------|----------------------|---------------------|----------------------|
| | G = Conservatives | | | G = Liberals | | |
| | GNH | ANES | Pooled | GNH | ANES | Pooled |
| True Mean Conservatives | 1.02**** (0.097) | 0.98**** (0.133) | 0.96**** (0.076) | -0.21**** (0.060) | -0.19 (0.116) | -0.25**** (0.060) |
| True Mean Liberals | -0.35*** (0.106) | -0.86**** (0.134) | -0.58**** (0.131) | 0.987**** (0.066) | 0.39*** (0.106) | 0.73**** (0.135) |
| Constant | 1.51*** (0.195) | 3.35**** (0.269) | 2.35**** (0.279) | 0.69**** (0.122) | 2.58**** (0.249) | 1.56**** (0.270) |
| R-squared | 0.83 | 0.53 | 0.66 | 0.92 | 0.32 | 0.68 |
| Obs. (Clusters) | 45 (45) | 66 (10) | 111 (55) | 45 (45) | 66 (10) | 111 (55) |

Notes: Std. errors in parentheses, clustered at the issue level. *, **, ***, and **** denote significance at the 10% level, 5%, 1%, and 0.1% level, respectively. In pooled specifications, we include a dummy variable indicating whether the observation came from ANES data set.

Table 2: Average Representativeness of Tail Positions Predicts Beliefs

| OLS Predicting Believed Mean of Group G | | | | | | |
|---|--------------------|--------------------|--------------------|--------------------|--------------------|---------------------|
| | G = Conservatives | | | G = Liberals | | |
| | GNH | ANES | Pooled | GNH | ANES | Pooled |
| True Mean of G | 0.78**** (0.06) | 0.24** (0.08) | 0.51**** (0.09) | 0.72**** (0.05) | 0.18*** (0.05) | 0.41**** (0.10) |
| Θ | 0.19** (0.07) | 0.55** (0.22) | 0.25*** (0.08) | -0.14** (0.06) | -0.12* (0.07) | -0.24**** (0.05) |
| Constant | 1.01**** (0.26) | 2.60**** (0.45) | 1.84**** (0.29) | 0.93**** (0.22) | 2.71**** (0.25) | 2.01**** (0.32) |
| R-squared | 0.82 | 0.48 | 0.60 | 0.91 | 0.31 | 0.70 |
| Obs. (Clusters) | 45 (45) | 66 (10) | 111 (55) | 45 (45) | 66 (10) | 111 (55) |

Notes: Std. errors in parentheses, clustered at the issue level. *, **, ***, and **** denote significance at the 10% level, 5%, 1%, and 0.1% level, respectively. In pooled specifications, we include a dummy variable indicating whether the observation came from ANES data set.

Figure 5 summarizes our evaluation of the two models for the GNH data, using the mean squared prediction error (MSPE) as a measure of the magnitude of errors.³⁰ It illustrates two findings. First, the representativeness model with $d = 4$ or $d = 5$ produces smaller MSPE than the accurate beliefs benchmark ($d = 6$). Second, the representativeness model also outperforms the likelihood model for these levels of d . In fact, while the likelihood model produces smaller MSPE than the stereotype model for small values of d , for each group, the best representativeness-based model produces smaller MSPE errors than the best likelihood-based model. And, while the best representativeness-based model is a better predictor of observed beliefs than the accurate beliefs benchmark for both liberals and conservatives in terms of MSPE, the best likelihood-based model never beats the accurate beliefs benchmark in terms of MSPE. Critically, the superior performance of the representativeness model comes precisely from its ability to capture the fact that beliefs are slanted toward representative, extreme group types.

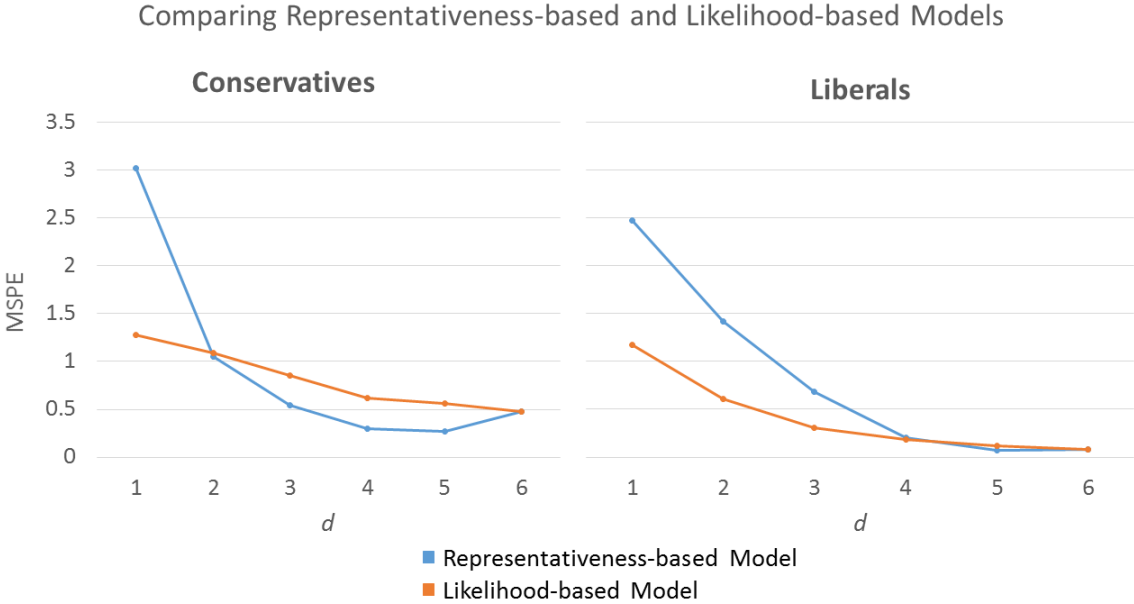


Figure 5: Comparing the Representativeness-Based and Likelihood-Based Models

We present several additional results in the online Appendix H, showing that similar

³⁰To compute MSPE for a given group and a given value of d , we subtract the model prediction for each observation from the observed data, square this difference, then take the mean of these squared differences.

results hold for the ANES data set (see Tables 9 and 10). We also explore other measures of prediction errors (e.g., mean prediction error, see Tables 7, 8, 9, and 10).

5 Conclusion

We presented a model of stereotypical thinking, in which decision makers making predictions about a group overweight the group’s most distinctive types. These overweighted types are not the most likely ones given the DM’s data, but rather the most representative ones, in the sense of being the most diagnostic of the group relative to other groups. Representativeness implies that what is most distinctive of a group depends on what group it is compared to. We presented experimental evidence that confirms this context dependence in recall-based assessments of groups. Finally, we evaluated the predictions of the model using political data from existing large scale surveys. We find context-dependence to be a key feature of beliefs. Given the richness of the political data, we can go a step further and identify a role for representativeness in particular. As the representativeness of tail types increases, beliefs of a group are distorted in the direction of that tail. A truncation model where the decision-maker neglects least representative types in forming beliefs about a group fits the data better than the accurate beliefs benchmark.

Our approach provides a parsimonious and psychologically founded account of how DMs generate simplified representations of reality, from social groups to stock returns, and offers a unified account of disparate pieces of evidence relating to this type of uncertainty. The model captures the central fact that stereotypes highlight the greatest difference between groups, thus explaining why some stereotypes are very accurate, while others lack validity. Still, stereotypes often have a “kernel of truth”, when they are based on systematic – even if small – differences between groups.

This same logic allows us to describe a number of heuristics and psychological biases, many of which arise in the context of prediction problems. Our model generates both base-rate neglect and confirmation bias (and makes novel predictions for when they occur). To our knowledge, ours is the first model to reconcile these two patterns of behavior, and in fact shows they both arise out of the assumption of representativeness-based recall.

Our model is based on representativeness and does not capture all the features of stereotypical thinking. However, it captures perhaps the central feature: when we think of a group, we focus on what is most distinctive about it, and neglect or underweight the rest.

References:

- Adorno, Theodor, Else Frenkel-Brunswik, Daniel Levinson, and Nevitt Sanford. 1950. *The Authoritarian Personality*. New York, NY: Harper & Row.
- Arrow, Kenneth. 1973. The Theory of Discrimination. In O. Ashenfelter and A. Rees, eds. *Discrimination in Labor Markets*. Princeton, N.J.: Princeton University Press: 3 – 33.
- Barberis, Nicholas, Andrei Shleifer, and Robert Vishny. 1998. “A Model of Investor Sentiment.” *Journal of Financial Economics* 49 (3): 307 – 343.
- Benjamin, Dan, James Choi, and Joshua Strickland. 2010. "Social Identity and Preferences." *American Economic Review* 100 (4): 1913 – 1928.
- Benjamin, Dan, Matthew Rabin, and Collin Raymond. 2015. "A Model of Non-Belief in the Law of Large Numbers." *Journal of the European Economic Association*, forthcoming.
- Bertrand, Marianne. 2011. “New Perspectives on Gender” in O. Ashenfelter and D. Card eds, *Handbook of Labor Economics*, 4 (B): 1543 – 1590.
- Bodoh-Creed, Aaron, Dan Benjamin, and Matthew Rabin. 2013. “The Dynamics of Base-Rate Neglect.” Mimeo Haas Business School.
- Bohnet, Iris, Alexandra van Geen, and Max Bazerman. 2015. “When Performance Trumps Gender Bias: Joint Versus Separate Evaluation.” *Management Science*, forthcoming.
- Bordalo, Pedro, Nicola Gennaioli, and Andrei Shleifer. 2012. ”Salience Theory of Choice under Risk.” *Quarterly Journal of Economics* 127 (3): 1243 – 1285.
- Bordalo, Pedro, Nicola Gennaioli, and Andrei Shleifer. 2013. “Salience and Consumer Choice.” *Journal of Political Economy* 121 (5): 803 – 843.
- Buser, Thomas, Muriel Niederle, and Hessel Oosterbeek. 2014. “Gender, Competitiveness and Career Choices.” *Quarterly Journal of Economics* 129 (3): 1409 – 1447.
- Carrell, Scott, Marianne Page, and James West. 2010. “Sex and Science: How Professor Gender Perpetuates the Gender Gap.” *Quarterly Journal of Economics* 125 (3): 1101 – 1144.
- Chan, Wayne et al. 2012. “Stereotypes of Age Differences in Personality Traits: Universal and Accurate?” *Journal of Personality and Social Psychology* 103(6): 1050 – 1066.

- Chen, Yan, Sherry Xin Liu, Tracy Xiao, and Margaret Shih. 2014. “Which Hat to Wear? Impact of Natural Identities on Coordination and Cooperation.” *Games and Economic Behavior* 84: 58 – 86.
- Cunningham, Tom. 2013. “Comparisons and Choice.” Unpublished manuscript, Stockholm University.
- Fryer, Roland, and Matthew Jackson. 2008. “A Categorical Model of Cognition and Biased Decision-Making.” *B.E. Journal of Theoretical Economics* 8(1): 1 – 42.
- Fryer, Roland, and Steven Levitt. 2010. “An Empirical Analysis of the Gender Gap in Mathematics.” *American Economic Journal, Applied Economics* 2(2): 210 – 240.
- Gennaioli, Nicola, and Andrei Shleifer. 2010. “What Comes to Mind.” *Quarterly Journal of Economics* 125 (4): 1399 – 1433.
- Gennaioli, Nicola, Andrei Shleifer, and Robert Vishny. 2012. “Neglected Risks, Financial Innovation, and Financial Fragility.” *Journal of Financial Economics* 104(3): 452 – 468.
- Gennaioli, Nicola, Andrei Shleifer, and Robert Vishny. 2015. “Neglected Risks: The Psychology of Financial Crises.” *American Economic Review, Papers & Proceedings* 105 (5): 310 – 14.
- Glover Dyland, Pallais Amanda, and Pariente William. 2015. “Discrimination as a Self-Fulfilling Prophecy: Evidence from French Grocery Stores.” Working paper.
- Goldin, Claudia, Lawrence Katz, and Ilyana Kuziemko. 2006. “The Homecoming of American College Women: The Reversal of the College Gender Gap.” *Journal of Economic Perspectives* 20 (4): 133 – 156.
- Graham, Jesse, Brian Nosek, and Jonathan Haidt. 2012. “The Moral Stereotypes of Liberals and Conservatives: Exaggeration of Differences across the Political Spectrum.” *PLOS One* 7 (12): 1 – 13.
- Grether, David. 1980. “Bayes Rule as a Descriptive Model: The Representativeness Heuristic.” *Quarterly Journal of Economics* 95 (3): 537 – 557.
- Griffin, Dale, and Amos Tversky. 1992. “The Weighing of Evidence and the Determinants of Confidence.” *Cognitive Psychology* 24 (3): 411 – 435.

- Guiso, Luigi, Ferdinando Monte, Paola Sapienza, and Luigi Zingales. 2008. "Culture, Gender, and Math." *Science* 320 (5880): 1164 – 1165.
- Hilton, James, and William Von Hippel. 1996. "Stereotypes." *Annual Review of Psychology* 47 (1): 237 – 271.
- Hyde, Janet, Sara Lindberg, Marcia Linn, Amy Ellis, and Caroline Williams. 2008. "Gender Similarities Characterize Math Performance." *Science* 321 (5888): 494 – 495.
- Judd, Charles, and Bernardette Park. 1993. "Definition and Assessment of Accuracy in Social Stereotypes." *Psychological Review* 100 (1): 109 – 128.
- Jussim, Lee, Jarret Crawford, Stephanie Anglin, John Chambers, Sean Stevens, and Florette Cohen. 2015. "Stereotype Accuracy: One of the Largest and Most Replicable Effects in All of Social Psychology" in *The Handbook of Prejudice, Stereotyping, and Discrimination*, Todd Nelson, editor. Mahwah, NJ: Lawrence Erlbaum Publishing.
- Kahneman, Daniel, and Amos Tversky. 1972. "Subjective Probability: A Judgment of Representativeness." *Cognitive Psychology* 3 (3): 430 – 454.
- Kahneman, Daniel, and Amos Tversky. 1973. "On the Psychology of Prediction." *Psychological Review* 80 (4): 237 – 251.
- Kersten, Daniel, Pascal Mamassian, and Alan Yuille. 2004. "Object Perception as Bayesian Inference." *Annual Review of Psychology*, 55: 271 – 304.
- Klein, Richard et al. 2014. "Investigating Variation in Replicability: A Many Labs Replication Project" *Social Psychology*, 45 (3): 142 – 152.
- Lippmann, Walter. 1922. *Public Opinion*. New York, NY: Harcourt.
- Lundberg, Shelly, and Richard Startz. 1983. "Private Discrimination and Social Intervention in Competitive Labor Markets." *American Economic Review* 73 (3): 340-347.
- Madon, Stephanie, Max Guyll, Kathy Aboufadel, Eulices Montiel, Alison Smith, Polly Palumbo, and Lee Jussim. 2001. "Ethnic and National Stereotypes: The Princeton Trilogy Revisited and Revised." *Personality and Social Psychological Bulletin* 27 (8): 996 – 1010.

- Mullainathan, Sendhil. 2002. "Thinking through Categories.", Working Paper, Harvard University.
- Ortoleva, Pietro, and Erik Snowberg. 2015. "Overconfidence in Political Behavior." *American Economic Review* 105 (2): 504 – 535.
- Phelps, Edmund. 1972. "The Statistical Theory of Racism and Sexism." *American Economic Review* 62 (4): 659 – 661.
- Pope, Devin, and Justin Sydnor. 2010. "Geographic Variation in the Gender Differences in Test Scores." *Journal of Economic Perspectives* 24(2): 95 – 108.
- Rabin, Matthew. 2002. "Inference by Believers in the Law of Small Numbers," *Quarterly Journal of Economics* 117 (3): 775 – 816.
- Rabin, Matthew, and Joel Schrag. 1999. "First Impressions Matter: A Model of Confirmatory Bias." *Quarterly Journal of Economics* 114 (1): 37 – 82.
- Rabin, Matthew, and Dimitri Vayanos. 2010. "The Gambler's and Hot-Hand Fallacies: Theory and Applications" *Review of Economic Studies* 77 (2): 730 – 778.
- Reuben, Ernesto, Paola Sapienza, and Luigi Zingales. 2014. "How Stereotypes Impair Women's Careers in Science." *Proceedings of the National Academy of Sciences* 111 (12): 4403 – 4408.
- Rosch, Eleanor. 1973. "Natural Categories." *Cognitive Psychology* 4 (3): 328 – 350.
- Schneider, David. 2004. *The Psychology of Stereotyping*. New York, NY: The Guilford Press.
- Schneider, David, Albert Hastorf, and Phoebe Ellsworth. 1979. *Person Perception* (2nd ed.). Reading, MA: Addison-Wesley.
- Shih, Margaret, Todd Pittinsky, and Nalini Ambady. 1999. "Stereotype Susceptibility: Identity Salience and Shifts in Quantitative Performance." *Psychological Science* 10 (1): 80 – 83.
- Shih, Margaret, Todd Pittinsky, and Amy Trahan. 2006. "Domain-specific Effects of Stereotypes on Performance." *Self and Identity* 5 (1): 1 – 14.

- Steele, Claude. 2010. *Whistling Vivaldi: How Stereotypes Affect Us and What We Can Do*. New York, NY: W. W. Norton & Company.
- Tenenbaum, Joshua, and Thomas Griffiths. 2001. "The Rational Basis of Representativeness." *23rd Annual Conference of the Cognitive Science Society*, 1036 – 1041.
- Tversky, Amos, and Daniel Kahneman. 1983. "Extensional versus Intuitive Reasoning: the Conjunction Fallacy in Probability Judgment." *Psychological Review* 90 (4): 293 – 315.
- Weinberger, Catherin. 2005. "Is the Science and Engineering Workforce Drawn from the Far Upper Tail of the Math Ability Distribution?" Working Paper, UCSB.
- Westfall, Jacob, Leaf Van Boven, John Chambers, and Charles Judd. 2015. "Perceiving Political Polarization in the United States: Party Identity Strength and Attitude Extremity Exacerbate the Perceived Partisan Divide." *Psychological Science* 10 (2): 145 – 158.

A Proofs

Proposition 1. By Definition 1, the representativeness ranking of types for G is the opposite of that for $-G$. Thus, if $t^{max,G} = \operatorname{argmax}_t \frac{\pi_{t,G}}{\pi_{t,-G}}$ is the most representative type for G , then it is also the least representative type for $-G$.

Suppose now that $\pi_{t,G} > \pi_{t',G}$ if and only if $\pi_{t,-G} > \pi_{t',-G}$ (case i)). Then, both groups share the same modal type t_{mod} . Because $\pi_G \neq \pi_{-G}$, it follows that not all types are equally representative. Because the representativeness ranking is opposite for the two groups, t_{mod} can coincide with the most representative type for at most one of the groups.

Consider now the case where $\pi_{t,G} > \pi_{t',G}$ if and only if $\pi_{t,-G} < \pi_{t',-G}$ (case ii)). Then, it also follows that $\pi_{t,G} > \pi_{t',G}$ if and only if $\pi_{t,G}/\pi_{t,-G} > \pi_{t',G}/\pi_{t',-G}$ so that likelihood and representativeness rankings coincide for each group. In particular, the most representative type coincides with the modal type for each group. ■

Proposition 2. Index the types $t \in \{1, \dots, T\}$ according to the underlying cardinal relation. Suppose first the likelihood ratio $\pi_{t,G}/\pi_{t,-G}$ is monotonically decreasing in t . Then it follows that $\pi_{t,-G}$ first order stochastically dominates $\pi_{t,G}$, so that $\mathbb{E}(t|G)$ is lower than $\mathbb{E}(t|-G)$, and therefore lower than the unconditional mean, $\mathbb{E}(t|G) < \mathbb{E}(t)$ (recall that $\mathbb{E}(t) = \mathbb{E}(t|\Omega)$ where $\Omega = G \cup -G$). Moreover, the ordering of types by representativeness coincides with the cardinal ordering of types, so that, for any function h_t , we have

$$h_t(\pi_G/\pi_{-G}) \leq h_{t'}(\pi_G/\pi_{-G}) \text{ iff } t > t'$$

where the first inequality is strict for at least some types. Consider now the likelihood ratio between the stereotypical distribution π_G^{st} and the undistorted distribution π_G :

$$\frac{\pi_{t,G}^{st}}{\pi_{t,G}} = \frac{h_t(\pi_G/\pi_{-G})}{\sum_{s \in T} \pi_{s,G} \cdot h_s(\pi_G/\pi_{-G})}$$

This likelihood ratio is (weakly) monotonically decreasing in t , implying that π_G f.o.s.d. π_G^{st} , and in particular that $\mathbb{E}^{st}(t|G) < \mathbb{E}(t|G)$.

If the the likelihood ratio is monotonically increasing in t , the same logic yields $\mathbb{E}(t|G) >$

$\mathbb{E}(t)$. Moreover, the ordering of types by representativeness coincides with the inverse of the cardinal ordering of types, so that now π_G^{st} f.o.s.d. π_G . It follows that $\mathbb{E}^{st}(t|G) > \mathbb{E}(t|G)$. ■

Proposition 3. Let the set of types be $T = \{0, \dots, T\}$ and consider for concreteness the case where T is even (the same proof goes through for T odd). Note first that the assumption that $\pi_{t,G}$ and $\pi_{t,-G}$ are symmetric around the midpoint $t = \frac{T}{2}$, namely that $\pi_{t,G} = \pi_{T-t,G}$, implies that representativeness $\pi_{t,G}/\pi_{t,-G}$ is also symmetric. This property ensures that the means $\mathbb{E}^{st}(t|G)$ and $\mathbb{E}^{st}(t|-G)$ are correctly estimated at $T/2$. Writing the weighting function as $f_t = h_t(\pi_G/\pi_{-G})/\sum_s \pi_{s,G}h_s(\pi_G/\pi_{-G})$ we find

$$\mathbb{E}^{st}(t|G) - \mathbb{E}(t|G) = \sum_t t \cdot \pi_{t,G} (f_t - 1) = \sum_{T-t} (T-t) \cdot \pi_{T-t,G} (f_{T-t} - 1) = \frac{T}{2} \sum_t \pi_{t,G} (f_t - 1) = 0$$

because $\mathbb{E}(f_t|G) = 1$. In contrast, variances are systematically distorted. To see this, write:

$$\text{Var}^{st}(t|G) - \text{Var}(t|G) = \sum_t \left(t - \frac{T}{2}\right)^2 \cdot \pi_{t,G} (f_t - 1) = \text{cov} \left(\left(t - \frac{T}{2}\right)^2, f_t - 1 \right)$$

If $\text{Var}(t|G) > \text{Var}(t|-G)$, so that representativeness is U-shaped, then $\text{cov} \left(\left(t - \frac{T}{2}\right)^2, f_t - 1 \right) > 0$ and variance gets exaggerated, $\text{Var}^{st}(t|G) > \text{Var}(t|G)$. Conversely, if $\text{Var}(t|G) < \text{Var}(t|-G)$ this implies $\text{Var}^{st}(t|G) < \text{Var}(t|G)$. ■

Proposition 4. 1) From Definition 2, we have

$$\mathbb{E}^{st}(t|G) = \mathbb{E}(t|G) - \sum_t t \cdot \pi_{t,G} \cdot (1 - f_t)$$

where $f_t = h_t(\pi_G/\pi_{-G})/\sum_s \pi_{s,G}h_s(\pi_G/\pi_{-G})$ is a differentiable weighing function. The factor $(1 - f_t)$ is increasing in the representativeness of type t for group $-G$, and equals zero when $\pi_{t,-G}/\pi_{t,G} = 1$. We will expand $(1 - f_t)$ to first order around $\pi_{t,-G}/\pi_{t,G} = 1$.

To do so, denote by $r_s = \frac{\pi_{s,G}}{\pi_{s,-G}}$ and by $x_s = 1/r_s$. The first order expansion around $\mathbf{x} = \mathbf{1}$ reads,

$$f_t(\mathbf{x}) \sim 1 + \sum_s \partial_{x_s} f_t|_{\mathbf{x}=\mathbf{1}} (x_s - 1)$$

We now characterise $\partial_{x_s} f_t$ around $\mathbf{x} = \mathbf{1}$. Denote $h_t(\mathbf{1}) = h$, and $\partial_{r_t} h_t|_{\mathbf{x}=\mathbf{1}} = \theta_1$ and $\partial_{r_s} h_t|_{\mathbf{x}=\mathbf{1}} = \theta_2$ for $s \neq t$, where by assumption $\theta_2 < 0 < \theta_1$. We have:

$$\begin{aligned} \partial_{x_t} f_t|_{\mathbf{x}=\mathbf{1}} &= - \left. \frac{\partial_{r_t} h_t \cdot (\sum_s \pi_{s,G} h_s) - h_t (\sum_s \pi_{s,G} \partial_{r_t} h_s)}{(\sum_s \pi_{s,G} h_s)^2} \right|_{\mathbf{x}=\mathbf{1}} \\ &= - \frac{\theta_1 h - h (\theta_1 \pi_{t,G} + \theta_2 (1 - \pi_{t,G}))}{h^2} \\ &= - \frac{(\theta_1 - \theta_2) (1 - \pi_{t,G})}{h} < 0 \end{aligned}$$

and

$$\partial_{x_s} f_t|_{\mathbf{x}=\mathbf{1}} = - \frac{(\theta_2 - \theta_1) \pi_{s,G}}{h} > 0$$

Inserting back into the first order expansion, we find

$$f_t(\mathbf{x}) \sim 1 - \frac{(\theta_1 - \theta_2)}{h} \cdot \left[(x_t - 1) - \sum_s \pi_{s,G} (x_s - 1) \right]$$

Note that, by construction, the average departure $x_s - 1$ is zero, namely $\sum_s \pi_{s,G} (x_s - 1) = \sum_s (\pi_{s,-G} - \pi_{s,-G}) = 0$. So we find $f_t(\mathbf{x}) \sim 1 - \frac{(\theta_1 - \theta_2)}{h} (x_t - 1)$. Plugging this approximation back into the expression for $\mathbb{E}^{st}(t|G)$ we find

$$\begin{aligned} \mathbb{E}^{st}(t|G) &\approx \mathbb{E}(t|G) - \sum_t t \cdot \pi_{t,G} \cdot \frac{(\theta_1 - \theta_2)}{h} (x_t - 1) \\ &= \mathbb{E}(t|G) \left(1 + \frac{(\theta_1 - \theta_2)}{h} \right) - \mathbb{E}(t|G) \left(\frac{\theta_1 - \theta_2}{h} \right) \end{aligned}$$

which yields the result, with $\kappa = \frac{\theta_1 - \theta_2}{h} > 0$.

2) Distinguish types into the right tail, $H = \{T - 2, \dots, T\}$, and the left tail, $L = \{1, \dots, T - 3\}$. For $X = L, H$, write $\pi_{X,G} = \sum_{t \in X} \pi_{t,G}$, $\bar{t}_X = \frac{\sum_{t \in X} t \cdot \pi_{t,G}}{\pi_{X,G}}$, and $f_X = \frac{\sum_{t \in X} \pi_{t,G} f_t}{\pi_{X,G}}$, so that $\pi_{L,G} f_L + \pi_{H,G} f_H = 1$ (we omit the group index from f_t for simplicity).

We then have

$$\mathbb{E}^{st}(t|G) = \mathbb{E}(t|G) + \sum_L t \cdot \pi_{t,G} \cdot (f_t - 1) + \sum_H t \cdot \pi_{t,G} \cdot (f_t - 1)$$

Under the approximation that f_t varies little within each tail, $f_t \approx f_X$ for $t \in X$, this becomes

$$\mathbb{E}^{st}(t|G) = \mathbb{E}(t|G) + (f_L - 1) \pi_L \bar{t}_L + (f_H - 1) \pi_H \bar{t}_H = \mathbb{E}(t|G) + (f_H - 1) \pi_H (\bar{t}_H - \bar{t}_L)$$

because $(f_L - 1) \pi_L = - (f_H - 1) \pi_H$.

Adapting the notation from point 1), we set $r_H = \frac{\sum_H \pi_{t,G}}{\sum_H \pi_{t,-G}}$ and $r_L = \frac{\sum_L \pi_{t,G}}{\sum_L \pi_{t,-G}}$. Expanding f_H around $r_H = 1$ we find

$$\begin{aligned} f_H(\mathbf{r}) &\approx 1 + \partial_{r_H} f_H|_{\mathbf{r}=\mathbf{1}} (r_H - 1) + \partial_{r_L} f_H|_{\mathbf{r}=\mathbf{1}} (r_L - 1) \\ &= 1 + (r_H - 1) \left[\partial_{r_H} f_H|_{\mathbf{r}=\mathbf{1}} - \partial_{r_L} f_H|_{\mathbf{r}=\mathbf{1}} \frac{\pi_{H,-G}}{\pi_{L,-G}} \right] \end{aligned}$$

Recall that by assumption $\partial_{r_H} f_H|_{\mathbf{r}=\mathbf{1}} > 0 > \partial_{r_L} f_H|_{\mathbf{r}=\mathbf{1}}$. Inserting this back into the approximate expression for $\mathbb{E}^{st}(t|G)$ we get

$$\mathbb{E}^{st}(t|G) \approx \mathbb{E}(t|G) + \lambda_G (r_H - 1)$$

where $\lambda_G = \pi_H (\bar{t}_H - \bar{t}_L) \left[\partial_{r_H} f_H|_{\mathbf{r}=\mathbf{1}} - \partial_{r_L} f_H|_{\mathbf{r}=\mathbf{1}} \frac{\pi_{H,-G}}{\pi_{L,-G}} \right]$ is positive. Thus, the believed mean $\mathbb{E}^{st}(t|G)$ increases in the representativeness r_H of the right tail for G , and is exaggerated when the right tail is more representative of G than the left tail. Replacing r_H with Θ yields equation (6). An analogous calculation, where we expand in $1/\Theta$, yields equation (7). ■

STEREOTYPES

SUPPLEMENTARY MATERIAL FOR ONLINE PUBLICATION

P. BORDALO, K. COFFMAN, N. GENNAIOLI, A. SHLEIFER

B Unordered Types

In many settings, decision makers must assess groups in terms of their distributions over unordered type spaces. For instance, one may be interested in the distribution of occupations, or of political views, or of beliefs of different social groups. Our model applies directly to these settings, provided the type space is specified, or at least implied, by the problem at hand. While there is no notion of “extreme” types in unordered type spaces, the central insight about how representativeness and likelihood combine to determine stereotype accuracy continues to hold: when groups are very similar, representative differences tend to be relatively unlikely, while when groups are different representative differences tend to be likely, and thus generate more accurate stereotypes.

To illustrate this logic in the context of unordered types, consider the formation of the stereotypes “Republicans are creationists” and “Democrats believe in Evolution”. In May 2012, Gallup conducted a public opinion poll assessing the beliefs about Evolution of members of the two main parties in the US. The results on the beliefs of Republicans and Democrats, largely unchanged in the three decades over which such polls have been conducted, are presented below:³¹

| | <i>Creationism</i> | <i>Evolution</i> | <i>Evolution guided by God</i> |
|--------------------|--------------------|------------------|--------------------------------|
| <i>Republicans</i> | 58% | 5% | 31% |
| <i>Democrats</i> | 41% | 19% | 32% |

The table shows that being a creationist is the distinguishing feature of the Republicans, not only because most Republicans are creationist but also because more Republicans are creationists than Democrats. In this sense, stereotyping a Republican as a creation-

³¹The three options were described as “God created Humans in present form in the last 10,000 years”, “Humans evolved, God has no part in process” and “Humans evolved, God guided the process”. See <http://www.gallup.com/poll/155003/Hold-Creationist-View-Human-Origins.aspx> for details.

ist yields a fairly accurate assessment. Formally, $t = \textit{Creationism}$ maximizes not only $\Pr(\textit{Republicans}|t)/\Pr(\textit{Democrats}|t)$ but also $\Pr(t|\textit{Republicans})$.

On the other hand, the distinguishing feature of the Democrats is to believe in the “standard” Darwinian Evolution of humans, a belief four times more prevalent than it is among Republicans. However, and perhaps surprisingly, only 19% of Democrats believe in Evolution. Most of them believe either in creationism (41%) or in Evolution guided by God (32%), just like Republicans do. Formally, $t = \textit{Evolution}$ maximizes $\Pr(\textit{Democrats}|t)/\Pr(\textit{Republicans}|t)$ but not $\Pr(t|\textit{Democrats})$. Evolution is not the most likely belief of Democrats, but rather the belief that occurs with the highest relative frequency. A stereotype-based prediction that a Democrat would believe in the standard evolutionary account of human origins, and would not believe in Creationism, is highly inaccurate.

Another example in this spirit is as follows. Suppose the DM must assess the time usage of Americans and Europeans. For the sake of simplicity, we consider only two types, namely $T = \{\text{time spent on work, time spent on vacation}\}$. The Americans work 49 weeks per year, so the conditional distribution of work versus vacation time is $\{0.94, 0.06\}$. In contrast, the Europeans work 47 weeks per year, with work habits $\{0.9, 0.1\}$. In both cases, work is by far the most likely activity. However, because the Americans’ work habits are more concentrated around their modal activity, the stereotypical American activity is work. Because Europeans have fatter vacation tails, their stereotypical activity is enjoying the dolce vita. This stereotype is inaccurate, precisely because the vast majority of time spent by Europeans is at work. Still, due to its higher representativeness, vacationing is the distinctive mark of Europeans, which renders the image of holidays highly available when thinking of that group.

C Multidimensional Types

In the real world, the types describing a group are multidimensional. Members of social groups vary in their occupation, education and income. Firms differ in their sector, location and management style. While in some cases only one dimension is relevant for the judgment at hand, in other cases multiple dimensions need to be considered. In these judgments, forming an appropriate model requires DM's to properly weigh the different dimensions. Representativeness has significant implications for this process. In particular, in many cases, the “kernel of truth” logic carries through to the case of multiple dimensions. Stereotypes are formed along the dimensions in which the groups differ most, although the DM focuses on proportional differences rather than absolute differences. As in the unidimensional case, stereotypes are context dependent in the sense that the dimensions along which a group is stereotyped depends on the other group it is compared to.

We focus on the special case in which there are two dimensions. A type consists of a vector (t_1, t_2) of two dimensions, where $t_i \in T_i$ for $i = 1, 2$. Denote by $\pi_{(t_1, t_2), G}$ and $\pi_{(t_1, t_2), -G}$ the joint probability densities in groups G and $-G$, respectively, which are defined over the set of types $T = T_1 \times T_2$. The representativeness of (t_1, t_2) for group G is given by:

$$R_G(t_1, t_2) \equiv \frac{\pi_{(t_1, t_2), G}}{\pi_{(t_1, t_2), -G}} = \frac{\pi_{t_1, G}}{\pi_{t_1, -G}} \cdot \frac{\pi_{t_2, (G, t_1)}}{\pi_{t_2, (-G, t_1)}}. \quad (8)$$

where $\pi_{t_2, (G, t_1)} = \Pr(t_2 | G, t_1)$. In light of Equation (8), then, we can immediately observe:

Lemma 1 *Suppose that $d < |T_1| \times |T_2|$ and that $\pi_{t_1, G} \neq \pi_{t_1, -G}$ for some $t_1 \in T_1$.*

i) If $\pi_{t_2, (G, t_1)} = \pi_{t_2, (-G, t_1)}$ for all t_1 and t_2 , then the stereotype for group G selects a subset of values for t_1 while allowing for all possible values of t_2 .

ii) If instead $\pi_{t_2, (G, t_1)} \neq \pi_{t_2, (-G, t_1)}$ for some t_1 and t_2 , then the stereotype for group G selects a subset of the most representative values of t_1 and t_2 .

Proof. If $\pi_{t_2, (G, t_1)} = \pi_{t_2, (-G, t_1)}$ for all t_1 and t_2 , as in case i), it follows from Equation 8 that $R_G(t_1, t_2) = R_G(t_1)$ (and similarly, $R_{-G}(t_1, t_2) = R_{-G}(t_1)$) for all t_1, t_2 . However, because $\pi_{t_1, G} \neq \pi_{t_1, -G}$ for some t_1 , it must be that $R_G(t_1) > R_G(t'_1)$ for some t_1, t'_1 . As a consequence, for d sufficiently small, the stereotype of G consists of a truncation $T_1^{st} \times T_2$,

where T_1^{st} includes only the types t_1 that have sufficiently high $R_G(t_1)$. The type space T_2 is not truncated because ties are included in the stereotype.

If instead $\pi_{t_2,(G,t_1)} \neq \pi_{t_2,(-G,t_1)}$ for some t_1 and t_2 , Equation 8 implies a strict representativeness ranking in at least a subset of types in $\{t_1\} \times T_2$. Thus, there exists $d < N$ such that some type in $\{t_1\} \times T_2$ is truncated and others are not. Similarly, because $\pi_{t_1,G} \neq \pi_{t_1,-G}$ for some t_1 , for given d some types in T_1 are truncated. Together, these observations imply that the stereotype for G generically implies truncations along both dimensions. ■

This result shows how the kernel of truth logic extends to multiple dimensions. When groups only differ along one dimension, namely when the distribution of t_2 is identical across groups conditional on t_1 (case i), the stereotype is formed along that dimension, in the sense that it highlights group differences in t_1 only. Suppose for instance that t_1 indexes education while t_2 indexes welfare status. If all groups are equally likely to be on welfare conditional on education, stereotypes exaggerate educational differences but the welfare status is correctly represented (conditional on education types that come to mind).³²

When instead groups differ along both dimensions (case ii), stereotypes highlight differences along both dimensions. In the context of the previous example, if the less educated group is *also* conditionally more likely to be on welfare, then it is stereotyped as “uneducated and on welfare”, while the other group is stereotyped as “educated and not on welfare”. Again, there is a kernel of truth in these stereotypes, but also an exaggeration of the correlation between education and being on welfare: people neglect that most elements of the less educated group are not on welfare, as well as the fact that a non-trivial share of the more educated, and possibly larger, group are in fact on welfare.

Multidimensional stereotypes also raise new aspects of context dependence. Consider the stereotype of the red-haired Irish. This stereotype arises from comparing the Irish to a population (e.g., Europeans) with a much lower share of red haired people. Our model predicts that this stereotype should change when the Irish are compared to a group with a similar share of red-haired people, such as the Scots. When compared to the Scots, a more

³²Here the stereotype allows for all possible values of t_2 because of the tie breaking assumption in Definition 2. The result that in case *i*) stereotypes are not organized along t_2 would continue to hold under the alternative assumption of random tie breaking. Even in this case, in fact, there would be no systematic selection of values of t_2 in the stereotypes of different DMs.

plausible stereotype for the Irish is “Catholic” because religion is the dimension along which Irish and Scots differ the most.

Formally, suppose that groups are characterized by two dimensions: hair color (red r , other o), and religion (catholic c , other \hat{o}). The Irish have a share r_i of red haired people and a share c_i of catholics. Europeans have a share r_e of red haired people and a share c_e of catholics. Critically, the Irish have a much higher share of red haired people, $r_i > r_e$, while catholics are similarly prevalent along the two groups, namely $c_i = c_e$. Hair color and religion are statistically independent in both populations.

Consider the stereotypes formed by comparing the Irish to Europeans. Lemma 1 implies stereotypes depend on the joint distribution of these variables. Because $c_i = c_e$, the representativeness of different types for the Irish is then given by:

$$\begin{aligned} R_i(r, c) &= \frac{r_i \cdot c_i}{r_e \cdot c_e} = \frac{r_i}{r_e} = \frac{r_i \cdot (1 - c_i)}{r_e \cdot (1 - c_e)} = R_i(r, \hat{o}) > \\ &> R_i(o, c) = \frac{(1 - r_i) \cdot c_i}{(1 - r_e) \cdot c_e} = \frac{1 - r_i}{1 - r_e} = \frac{(1 - r_i) \cdot (1 - c_i)}{(1 - r_e) \cdot (1 - c_e)} = R_i(o, \hat{o}). \end{aligned}$$

The inequality follows because $r_i > r_e$ implies that $\frac{r_i}{r_e} > \frac{1 - r_i}{1 - r_e}$. As a consequence, when $d = 1$, the stereotype for the Irish contains the two equally representative types of (red haired, catholic) and (red haired, other). The stereotype differentiates the Irish from the Europeans along the color of hair dimension.

Suppose now that the Irish are compared to the Scots, who have a share r_s of red haired people and a share c_s of catholics. The Scots have a similar share of red haired people, $r_i = r_s$, while they have a much lower share of catholics, namely $c_i > c_s$. Consider the stereotype formed by comparing the Irish to the Scots. In this case, the representativeness of different types for the Irish is:

$$\begin{aligned} R_i(r, c) &= \frac{r_i \cdot c_i}{r_s \cdot c_s} = \frac{c_i}{c_s} = \frac{(1 - r_i) \cdot c_i}{(1 - r_s) \cdot c_s} = R_i(o, c) > \\ &> R_i(r, \hat{o}) = \frac{r_i \cdot (1 - c_i)}{r_s \cdot (1 - c_s)} = \frac{1 - c_i}{1 - c_s} = \frac{(1 - r_i) \cdot (1 - c_i)}{(1 - r_s) \cdot (1 - c_s)} = R_i(o, \hat{o}) \end{aligned}$$

Note that now $c_i > c_s$ implies that $\frac{c_i}{c_s} > \frac{1 - c_i}{1 - c_s}$. As a consequence, when $d = 1$, the stereotype for the Irish contains the two equally representative types of (red haired, catholic)

and (other, catholic). The dimensions along which the Irish stereotype is formed has changed: it differentiates the Irish from the Scots along the religion dimension, not along hair color.

In summary, because stereotypes are centered along the types for which the groups differ the most, the kernel of truth logic survives when types are multidimensional. The features that are perceived as characteristic of a group depend on the comparison group.

D Extension to Continuous Distributions

Many distributions of interest in economics can be usefully approximated by continuous probability distributions. Here we show how our results extend to this case. For simplicity, we only consider rank-based truncation, but the model is easily extended to smooth weighing.

D.1 Basic Setting

Let T be a continuous variable defined on the support $\bar{T} \subseteq R^k$. Denote by $t \in \bar{T}$ a realization of T which is distributed according to a density function $f(t) : \bar{T} \rightarrow R_+$. Denote by $f(t|G)$ and $f(t|-G)$, the distributions of t in G and $-G$, respectively. In line with Definition 1, we define representativeness as:

Definition 3 *The representativeness of $t \in \bar{T}$ for group G is measured by the ratio of the probability of G and $-G$ at $T = t$, where $-G = \Omega \setminus G$. Using Bayes' rule, this implies that representativeness increases on the likelihood ratio $f(t|G)/f(t|-G)$.*

In the continuous case, the exemplar for G is the realization t that is most informative about G . For one dimensional variables, the exemplar for G is $\sup(\bar{T})$ if the likelihood ratio is monotone increasing, or $\inf(\bar{T})$ if the likelihood ratio is monotone decreasing, just as in Proposition 2.

The DM constructs the stereotype by recalling the most representative values of t until the recalled probability mass is equal to the bounded memory parameter $\delta \in [0, 1]$. When $\delta = 0$, the DM only recalls the most representative type. When $\delta = 1$ the DM recalls the entire support \bar{T} and his beliefs are correct. When δ is between 0 and 1, we are in an intermediate case.

Definition 4 *Given a group G and a threshold $c \in R$, define the set $\bar{T}_G(c) = \left\{ t \in \bar{T} \mid \frac{f(t|G)}{f(t|-G)} \geq c \right\}$. The DM forms his beliefs using a truncated distribution in $\bar{T}_G(c(\delta))$ where $c(\delta)$ solves:*

$$\int_{t \in \bar{T}_G(c(\delta))} f(t|G) dt = \delta.$$

The logic is similar to that of Definition 2, with the only difference that now the memory constraint acts on the recalled probability mass and not on the measure of states, which would be problematic to compute when distributions have unbounded support. This feature yields and additional (and potentially testable) prediction that changes in the distribution typically change also the support of the stereotype by triggering the DM to recall or forget some states, even when the states' relative representativeness does not change.

D.2 The Normal Case

When $f(t|G)$ and $f(t|-G)$ are univariate normal, with means μ_G , μ_{-G} and variances σ_G , σ_{-G} , the stereotype of G is easy to characterize.

Proposition 5 *In the normal case, the stereotype works as follows:*

i) Suppose $\sigma_G = \sigma_{-G} = \sigma$. Then, if $\mu_G > \mu_{-G}$ the stereotype for G is $\bar{T}_G = [t_G, +\infty)$, where t_G decreases with δ . Moreover, $\mathbb{E}^{st}(t|G) > \mu_G > \mu_{-G} > \mathbb{E}^{st}(t|-G)$.

If instead $\mu_G < \mu_{-G}$, the stereotype for G is $\bar{T}_G = (-\infty, t_G]$, where t_G now increases with δ . Moreover, $\mathbb{E}^{st}(t|G) < \mu_G < \mu_{-G} < \mathbb{E}^{st}(t|-G)$. In both cases, $Var^{st}(t|G) < Var(t|G)$ and $Var^{st}(t|-G) < Var(t|-G)$.

ii) Suppose that $\sigma_G < \sigma_{-G}$. Then, the stereotype for G is $\bar{T}_G = [\underline{t}_G, \bar{t}_G]$ where \underline{t}_G decreases and \bar{t}_G increases with δ . Moreover, $Var^{st}(t|G) < Var(t|G)$.

iii) Suppose that $\sigma_G > \sigma_{-G}$. Then, the stereotype for G is $\bar{T}_G = (-\infty, \underline{t}_G] \cup [\bar{t}_G, +\infty)$ where \underline{t}_G increases and \bar{t}_G decreases with δ . Moreover, $Var^{st}(t|G) > Var(t|G)$.

Proof. Let ρ_{μ, σ^2} denote the probability density of $N(\mu, \sigma^2)$, namely $\rho(t) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(t-\mu)^2}{2\sigma^2}}$. The exemplar \hat{t}_G of $G \equiv N(\mu_G, \sigma_G^2)$ relative to $-G \equiv N(\mu_{-G}, \sigma_{-G}^2)$ satisfies $\hat{t}_E = \operatorname{argmax}_t \frac{\rho_{\mu_G, \sigma_G^2}}{\rho_{\mu_{-G}, \sigma_{-G}^2}}$ where

$$\frac{\rho_{\mu_G, \sigma_G^2}}{\rho_{\mu_{-G}, \sigma_{-G}^2}} = \frac{\sigma_{-G}}{\sigma_G} \cdot \exp \left\{ -t^2 \left(\frac{1}{2\sigma_G^2} - \frac{1}{2\sigma_{-G}^2} \right) + t \left(\frac{\mu_G}{\sigma_G^2} - \frac{\mu_{-G}}{\sigma_{-G}^2} \right) - \left(\frac{\mu_G^2}{2\sigma_G^2} - \frac{\mu_{-G}^2}{2\sigma_{-G}^2} \right) \right\}$$

When $\sigma_G < \sigma_{-G}$, the function above has a single maximum in t , namely that which maximizes the parabola in the exponent, $\hat{t}_E = \frac{\frac{\mu_G}{\sigma_G^2} - \frac{\mu_{-G}}{\sigma_{-G}^2}}{\frac{1}{\sigma_G^2} - \frac{1}{\sigma_{-G}^2}}$ from which the result follows.

When $\sigma_G > \sigma_{-G}$, the function above is grows without bounds with $|t|$, so that $\hat{t}_G \in \{-\infty, +\infty\}$.

When $\sigma_G = \sigma_{-G} = \sigma$, the exemplar \hat{t}_G of $G \equiv N(\mu_G, \sigma^2)$ relative to $-G \equiv N(\mu_{-G}, \sigma^2)$ satisfies

$$\hat{t}_G = \operatorname{argmax}_t e^{-\frac{\mu_G^2 - \mu_{-G}^2}{2\sigma^2}} \cdot e^{\frac{t}{2\sigma^2}(\mu_G - \mu_{-G})}$$

so that $\hat{t}_G = -\infty$ if $\mu_G < \mu_{-G}$ and $\hat{t}_G = +\infty$ otherwise. If $\mu_G < \mu_{-G}$ all values of t are equally representative. ■

When the two distributions have the same variance, the stereotype is formed by truncating from the original distribution the least representative tail (as in Section 3). In fact, when the mean in G is above the mean in $-G$, the likelihood ratio is monotone increasing and the exemplar for G is $+\infty$; otherwise it is $-\infty$. In both cases, the exemplar is inaccurate because it relies on a highly representative but very low probability realization.

Figure 6, left panel, represents the distribution considered by the DM for the high mean group when traits are normally distributed with the same variance across groups.

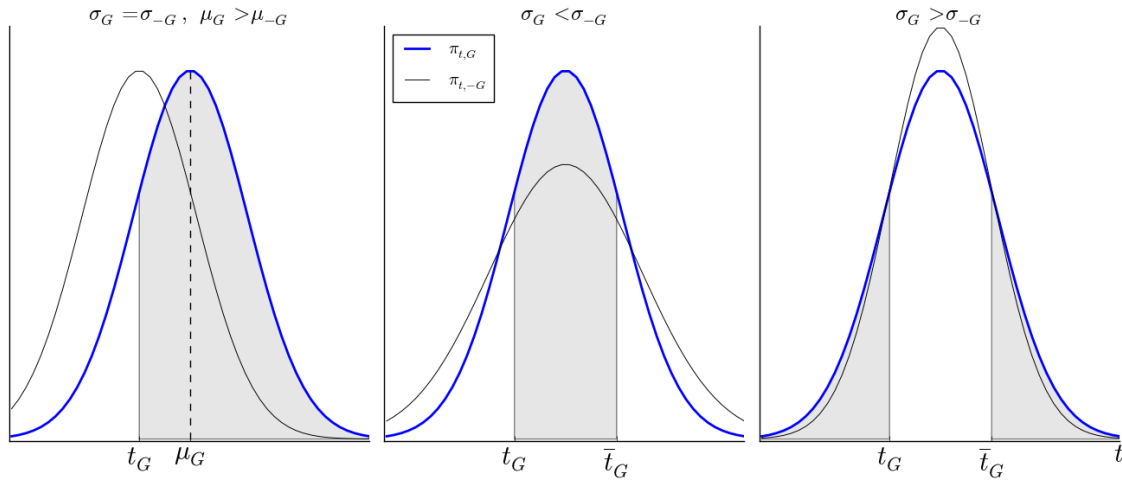


Figure 6: Stereotypes of a Normal distribution as a function of μ_{-G} and σ_{-G} .

Consider now case ii), where the variance of G is lower than that of $-G$, Figure 6, middle panel. The stereotype consists of an interval around an intermediate exemplar, denoted by \hat{t}_G . When the distribution in G is more concentrated than that in $-G$, the exemplar is accurate

and captures a relatively frequent, intermediate event. It is however somewhat distorted, because \hat{t}_G lies below the group's true mean μ_G if and only if $\mu_G < \mu_{-G}$. Interestingly, when the mean in the two groups is the same, the low variability group is represented by its correct mean, namely μ_G . Again, because the distinctive feature of group G is being more “average” than group $-G$, its stereotype neglects extreme elements and decreases within group variation.

Finally, consider case iii). Now the variance in G is higher than that in $-G$, Figure 6, right panel. As a consequence, both tails are exemplars and the stereotype includes both tails, truncating away an intermediate section of the distribution. This representation increases perceived volatility and thus captures the distinctive trait of G relative to $-G$, which is precisely its higher variability. Stereotyping now induces the DM to recall group G 's most extreme elements and to perceive G as more variable than it really is. This is a testable prediction of our model that stands in contrast with the previous cases, and with the common description that stereotypes reduce within-group variability (Hilton and Von Hippel 1996). However, it is consistent with the more basic intuition that stereotyping highlights the most distinctive features of group G , in this case its extreme elements. As an illustration of this mechanism, when thinking about stock returns, investors may think of positive scenarios where returns are high, or negative scenarios where returns are low, but neglect average returns, which are more typical of safer asset classes.

Consider now dynamic updating in this normal case. The DM receives information about the distributions $f(t|G)$ and $f(t|-G)$ over time. In each period k , a sample $(t_{G,k}, t_{-G,k})$ of outcomes is observed, drawn from the two groups. The history of observations up to period K is denoted by the vector $t^K = (t_{G,k}, t_{-G,k})_{k=1, \dots, K}$.

Based on t^K , and thus on the conditional distributions $f(t|W, t^K)$ for $W = G, -G$, the DM updates stereotypes and beliefs. In one tractable case, the $k = 0$ initial distribution $f(t|W)$ is also normal for $W = G, -G$. Formally, suppose that $t_W = \theta_W + \varepsilon_W$ where ε_W is i.i.d. normally distributed with mean 0 and variance v , and θ_W is the group specific mean. Initially, groups are believed to be identical, in the sense that both θ_G and θ_{-G} are normally distributed with mean 0 and variance γ . After observing $(t_{G,1}, t_{-G,1})$, the distribution of θ_W is updated according to Bayesian learning. Updating continues as progressively more

observations are learned. Thus, after observing the sample t^K , we have:

$$f(t|W, t^K) = \mathcal{N}\left(\frac{\gamma \cdot K}{v + \gamma \cdot K} \cdot \frac{\sum t_{W,k}}{K}; v \cdot \frac{v + \gamma \cdot (K + 1)}{v + \gamma \cdot K}\right). \quad (9)$$

The posterior mean for group W is an increasing function of the sample mean $\sum t_{W,k}/K$ for the same group. The variance of the posterior declines in sample size K , because the building of progressively more observations reduces the variance of θ_W , in turn reducing the variability of outcomes. However, and importantly, because the same number of observations is received for each group, both groups have the same variance in all periods.

Consider now how learning affects stereotypes. Proposition 5 implies:

Proposition 6 *At time K , the stereotype for group G is equal to $[t_G, +\infty)$ if $\sum t_{G,k} > \sum t_{-G,k}$ and to $(-\infty, t_G]$ if $\sum t_{G,k} < \sum t_{-G,k}$. As a result:*

i) Gradual improvement of the performance of group G does not improve that group's exemplar (and only marginally affects its stereotype) provided $\sum t_{G,k}$ stays below $\sum t_{-G,k}$. In particular, common improvements in the performance of G and $-G$ (which leave $\sum t_{G,k} - \sum t_{-G,k}$ constant) leave stereotypes unaffected.

ii) Small improvements in the relative performance of G that switch the sign of $\sum t_{G,k} - \sum t_{-G,k}$ have a drastic effect on stereotypes.

Proof. Since the variances of the sample populations G and $-G$ are equal, the stereotypes are fully determined by the sample means. From Proposition 5, if $\sum_t t_{G,k} > \sum_t t_{-G,k}$, then the sample mean of G is higher than that of $-G$, so that its exemplar is $\hat{t}_G = +\infty$. If instead $\sum_t t_{G,k} < \sum_t t_{-G,k}$, the exemplar of G is $\hat{t}_G = -\infty$. Cases i) and ii) follow directly from this. ■

Even in the normal case, stereotyping suffers from both under- and over-reaction to information. If new information does not change the ranking between group averages, exemplars do not change and stereotypes only respond marginally. Thus, even if a group gradually increases its average, its stereotype may remain very low. In contrast, even minor information can cause a strong over-reaction if it reverses the ranking between group averages.

E Likelihood, Availability, and Stereotypes

As we discussed in Section 2.2, our formulation of representativeness-based stereotypes leads in some instances to extreme predictions and, importantly, neglects other factors that influence what features come to mind when thinking about a group, such as likelihood and availability.³³ When stereotyping the occupation of a democratic voter, people think about “professor” rather than a “comparative literature professor.” While the latter is probably more representative, the former is more likely and thus comes to mind more easily.

In this section we show that our model can be easily adapted to account for some effects of likelihood on recall. When we do so, our predictions become less extreme, in the sense that stereotypes become centered around relatively more likely or available types, but the distortions of stereotypes still follow the logic of representativeness, as in our main analysis. This extension can also capture the effects of a crude measure of availability on recall. For simplicity, we focus on a rank-based truncation specification.

Suppose that the ease of recall of a type t for group G is given by:

$$R_k(t, G) = \frac{\pi_{t,G}}{\pi_{t,-G} + k} = \frac{1}{\frac{1}{R(t,G)} + k \cdot \frac{1}{\pi_{t,G}}} \quad (10)$$

where $k \geq 0$ and $R(t, G)$ is representativeness as defined in Definition 1. In Equation (10), the ease of recalling type t increases when that type is more representative, namely when $R(t, G)$ is higher, but also when type t is more likely in group G , namely when $\pi_{t,G}$ is higher. The value of k modulates the relative strength of these two effects: for small k , representativeness drives ease of recall, while for large k likelihood drives recall.³⁴

In this new formulation, the stereotype is formed as in Definition 2 except that now what comes to mind are the d types that are easiest to recall. When representative types are also likely, recall based on Equation (10) does not change the stereotype for group G . When instead representativeness and likelihood differ for group G , recall driven by $R_k(t, G)$ may

³³According to Kahneman and Frederick (2005) “the question of why thoughts become accessible – why particular ideas come to mind at particular times – has a long history in psychology and encompasses notions of stimulus salience, associative activation, selective attention, specific training, and priming”.

³⁴When $k = 0$, we are in a pure representativeness model. As k increases, likelihood becomes progressively more important in shaping recall relative to representativeness. As $k \rightarrow \infty$, only likelihood matters for shaping recall and stereotypes.

yield a different stereotype than a pure representativeness model.

To see how the model can capture some features of availability, note that the term $\pi_{t,G}$ in (10), and also in (2), may be broadly interpreted as capturing the availability, rather than just the frequency, of type t for group G . Formally, in the model of learning of Section F, we would assume that the estimate of $\pi_{t,G}$ is determined by the share of observations from G that are of type t , even if these observations are not independent. Thus, as the same episodes of terrorism are mentioned repeatedly in the news, their ease of recall is inflated. In this approach, availability is related to neglect of the correlation structure of information (as discussed in Section 2.2, the psychology of availability is beyond the scope of this paper).

The concrete implications of Equation (10) are best seen in the case where the type space is continuous, and more specifically when t is normally distributed in groups G and $-G$, with means μ_G, μ_{-G} respectively, and variance σ . In this case, the easiest to recall type t for group G is given by:

$$t_{E,G} = \operatorname{argmin}_t e^{\frac{(t-\mu_G)^2 - (t-\mu_{-G})^2}{2\sigma^2}} + k \cdot e^{\frac{(t-\mu_G)^2}{2\sigma^2}}$$

When $\mu_G > \mu_{-G}$, the easiest to recall type $t_{E,G}$ satisfies:

$$k \cdot (t_{E,G} - \mu_G) \cdot e^{\frac{(t_{E,G} - \mu_G)^2 + 2(\mu_G - \mu_{-G}) \cdot (t_{E,G} - \frac{\mu_G + \mu_{-G}}{2})}{2\sigma^2}} = \mu_G - \mu_{-G} \quad (11)$$

The left hand side of (11) is increasing in $t_{E,G}$, which implies that $t_{E,G}$ is a strictly increasing function of k satisfying $\lim_{k \rightarrow \infty} t_{E,G}(k) = \mu_G$ and $\lim_{k \rightarrow 0} t_{E,G}(k) = \infty$. In words, the group G with higher mean is stereotyped with an inflated assessment that goes in the direction of the most representative type $t = \infty$. The extent of this inflation increases as k gets smaller. The stereotype for group G in this case is an interval around the easiest to recall type that captures a total probability mass of δ (truncating both tails, but especially the left one). Moreover, as in the case $k = 0$, the stereotype has a lower variance than the true distribution. A corresponding result is obtained if group G has a lower mean than $-G$.

This analysis implies that the basic insights that stereotypes emphasise differences, and lead to base rate neglect, carry through to this case.³⁵

³⁵In the extended model given by (10), the parameters δ and k capture two natural types of bounds on recall: δ determines "how much" comes to mind (which might depend on effort), while k corresponds to the

F Stereotypes and Reaction to New Information

Stereotypes are hard to change, but they are far from immutable. For instance, stereotypes of immigrant populations change over time: in the early 20th century US, European Jews were stereotyped as religious and Asian immigrants were stereotyped as uneducated, yet both groups are stereotyped as high-achievers at the beginning of the 21st (Madon et. al., 2001). More recently, a rapid increase in the share of female doctors has coincided with shifting gender stereotypes in the medical profession. Medicine has historically been perceived as a stereotypical male profession, with women being viewed as less competent than their male counterparts (Decker 1986). However, this stereotype has faded, with specialties where women are more prevalent, such as pediatrics and dermatology, now being viewed as gender neutral (Couch and Sigler 2001). These patterns reflect at least in part changes in stereotypes in response to changes in reality. In fact, the experimental psychology literature documents that stereotypes change when individuals are faced with sufficiently pressing disconfirming information (Schneider 2004).

Our model can be naturally extended to investigate how stereotypes and beliefs change by the arrival of new information over time. To explore these dynamics, we suppose that at the outset, unlike in Section 2, the decision maker does not have perfect information about the categorical distribution $(\pi_{t,G})_{t=1,\dots,N}$ of the group G of interest, or about the distribution $(\pi_{t,-G})_{t=1,\dots,N}$ of the comparison group $-G$. Instead, the DM has priors over these distributions that are described by the Dirichlet distribution:

$$g[\pi_{t,W}, \alpha_{t,W}]_{t=t_1,\dots,t_N} = \frac{\Gamma(\sum_t \alpha_{t,W})}{\prod_t \Gamma(\alpha_{t,W})} \cdot \prod_t \pi_{t,W}^{\alpha_{t,W}-1}, \quad \text{for } W = G, -G,$$

which are conveniently conjugate to the categorical distributions assumed so far. Parameters $\alpha_G = (\alpha_{t,G})_{t=t_1,\dots,t_N}$ and $\alpha_{-G} = (\alpha_{t,-G})_{t=t_1,\dots,t_N}$ pin down the prior expectations of a Bayesian agent:

$$\Pr(T = t | \alpha_W) = \mathbb{E}(\pi_{t,W} | \alpha_W) = \frac{\alpha_{t,W}}{\sum_u \alpha_{u,W}}, \quad \text{for } W = G, -G. \quad (12)$$

In contrast to the Bayesian agent, the stereotype initially held by the DM depends on

relative weight of likelihood in recall, which may vary across people.

the probabilities in Equation (12) according to Definition 1. For simplicity, we set $\sum_t \alpha_{t,G} = \sum_t \alpha_{t,-G}$.

Suppose that a sample $n_W = (n_{1,W}, \dots, n_{N,W})$ is observed, where $n_{t,W}$ denotes the observation count in type t and let $\sum_t n_{t,W}$ be the total number of observations for group W . Then, a Bayesian's posterior probability of observing t is

$$\Pr(T = t | \alpha_W, n_W) = \mathbb{E}(\pi_{t,W} | \alpha_W, n_W) = \frac{\alpha_{t,W} + n_{t,W}}{\sum_u (\alpha_{u,W} + n_{u,W})}, \quad (13)$$

which is a weighted average of the prior probability of Equation (12) and the sample proportion $n_{t,W}/n_W$ of type t . As new observations arrive, the probability distribution in group W , and thus stereotypes, are updated according to Equation (13).³⁶

Consider how a DM influenced by representativeness updates beliefs. Again, we consider a rank-based truncation specification. Proposition 7 describes how new information changes the set of types that come to mind, shedding light on when and how stereotypes change. Proposition 8 considers the effect of information on probability assessments on types that are already included in the stereotype.

Proposition 7 *Suppose that the DM observes the same number of realizations from both groups, formally $\sum_u n_{u,G} = \sum_u n_{u,-G} = n$. Then:*

i) If for both groups all observations occur on the same type t that is initially non-representative for G , then this type does not become representative for G . Formally, if $n_{t,G} = n_{t,-G} = n$ for a type t such that $\alpha_{t,G}/\alpha_{t,-G} < 1$, then $\Pr(T = t | \alpha_W, n_G) / \Pr(T = t | \alpha_W, n_{-G}) < 1$ for all n .

ii) If all observations for G occur in a non representative type for G , while those for $-G$ occur in a type that is representative for G , then for a sufficiently large number of observations the stereotype for G changes. Formally, if $n_{t,G} = n$ for a type t such that $\alpha_{t,G}/\alpha_{t,-G} < 1$, while $n_{t',-G} = n$ for a type t' such that $\alpha_{t',G}/\alpha_{t',-G} > 1$, then for n sufficiently

³⁶While we assume for simplicity that updating is Bayesian, the representativeness mechanism that links priors to stereotypes can naturally be coupled with a non-Bayesian updating process. Psychologists have documented a tendency to search for information that confirms one's beliefs (Lord, Ross and Lepper 1979, Nickerson 1998). Schwartzstein (2014) proposes a model of biased learning in which information is used to update beliefs only about dimensions that are attended to.

large $\Pr(T = t'|\alpha_W, n_G)/\Pr(T = t'|\alpha_W, n_{-G}) < 1 < \Pr(T = t|\alpha_W, n_G)/\Pr(T = t|\alpha_W, n_{-G})$.

Proof. We assume that the same number of observations are received at each stage of the learning process for both groups G and $-G$. This assumption is not restrictive, since only the relative frequency of observations matters. In particular, all probabilities remain unchanged if the sample size of one group is scaled up relative to the sample size of the other. Thus we can set $\sum_{t'} a_{t',G} = \sum_{t'} a_{t',-G} = a$ and $\sum_{t'} n_{t',G} = \sum_{t'} n_{t',-G} = n$.

Representativeness of a type t is now measured by the ratio

$$\frac{\Pr(X = x|\alpha_S, n_S)}{\Pr(X = x|\alpha_{-G}, n_{-G})} = \frac{\alpha_{t,G} + n_{t,G}}{\alpha_{t,-G} + n_{t,-G}}$$

where $\alpha_S = (\alpha_{t,S})_{t \in T}$ are the priors for group S and $n_S = (n_{t,S})_{t \in T}$ is the sample for group S .

Consider case i) where all observations occur in type t , so that $n_{t,G} = n$ and $n_{t',G} = 0$ for $t' \neq t$, and similarly for $-G$. Then the representativeness of types other than t does not change, while the representativeness of t is $(\alpha_{t,G} + n)/(\alpha_{t,-G} + n_{t,-G})$. This tends to one monotonically as n increases. Therefore, if $a_{t,G}/a_{t,-G} < 1$ then $(a_{t,G} + n)/(a_{t,-G} + n) < 1$ for all n : namely, if t is non-representative to begin with, then no amount of observations of t in population G (when accompanied by observations of t in population $-G$) will make t representative for G .

Consider now case ii), where all observations in G occur in a non-representative type t while all observations in $-G$ occur in a representative (for G) type t' . In that case, the representativeness of t for group G increases as $(a_{t,G} + n)/(a_{t,-G})$, while the representativeness of t' for group G decreases as $(a_{t',G} + n)/(a_{t',-G} + n)$. The result follows. ■

The stereotype for a group does not necessarily change if the new observations are contrary to the initial stereotype. The stereotype is modified only if the new information for G and $-G$ is sufficiently different, as in case ii). Only when a disproportionate number of non-stereotypical observations occur for group G do these previously neglected types become sufficiently more likely in, and thus representative of, G .

Proposition 7 describes how much contrary data is needed in order to change a stereotype. For example, a process of economic development can significantly improve the livelihoods of all groups in a population, but it does not dispel the negative stereotype of a group that still includes a disproportional share of low income households. Instead, if a (subset of a) negatively stereotyped group outperforms the overall population, then its stereotype can change. Attitudes towards certain immigrant groups (Jews, Asians) have changed only as a subset of them overtook the overall population in terms of socioeconomic status. Attitudes towards women in the medical profession have changed only after a dramatic catch up in the number of female doctors, particularly in specialties in which female doctors are as frequent as male doctors (Noori and Weseley 2011).

We now consider how the initial stereotype for group G (formally, the priors over G and $-G$) affects the way in which the DM processes new information about G . We assume the support of the stereotype for G is fixed, and explore how the DM reacts to further information about G (note that, once the support of G 's stereotype is determined, information about $-G$ plays no role in determining the stereotypical distribution).

Proposition 8 *Let $d > 1$. Suppose that one observation about type t is received in group G (formally, $n = n_{t,G} = 1$). Then:*

i) If t belongs to the stereotype of G and its probability is sufficiently low, the DM over-reacts (relative to the Bayesian) in revising upward his assessment of t 's probability. Formally, there is a threshold $\nu \in (0, 1/2)$ such that the DM's assessed probability of t increases by more than under Bayesian updating if and only if $\alpha_{t,G} / \sum_u a_{u,G} < \nu$.

ii) If t does not belong to the stereotype of G , the DM does not update its probability at all, so he under-reacts relative to the Bayesian DM.

Proposition 8. Consider the case where a single observation of group G occurring in type t does not change the representativeness ranking of types – and thus the stereotype – for G .

If t is in the stereotype of G , then its estimated probability is $a_{t,G} / \sum_{t'=1}^d a_{t',G}$, which is boosted by a factor of $\sum_{t'=1}^N a_{t,G} / \sum_{t'=1}^d a_{t',G} > 1$, where d is the number of types in the stereotype. Suppose an observation occurs in type t . Its representativeness for G increases, and its assessed probability jumps to $(a_{t,G} + 1) / (\sum_{t'=1}^d a_{t',G} + 1)$. This corresponds to a larger

increase of assessed probability than that made by a Bayesian whenever

$$\frac{a_{t,G} + 1}{\sum_{t'=1}^d a_{t',G} + 1} - \frac{a_{t,G}}{\sum_{t'=1}^d a_{t',G}} > \frac{a_{t,G} + 1}{\sum_{t'=1}^N a_{t',G} + 1} - \frac{a_{t,G}}{\sum_{t'=1}^N a_{t',G}}$$

namely when

$$\frac{a_{t,G}}{\sum_{t'=1}^N a_{t',G}} < \frac{\sum_{t'=1}^d a_{t',G}}{1 + \sum_{t'=1}^d a_{t',G} + \sum_{t'=1}^N a_{t',G}} < \frac{1}{2}$$

The intuition is that the stereotype ignores some observations, it is as though the probability is being updated over a smaller sample size. Therefore, as long as the prior of t (in the stereotype) is not too large, the DM boosts it more than the Bayesian.

If t is not in the stereotype, then – given that the stereotype does not change – it does not become representative. Its assessed probability stays at zero, so the decision maker under-reacts to this observation relative to a Bayesian. ■

Proposition 8 indicates that stereotypes can both over and under-react to information. In case i), the DM strongly over-reacts to information confirming the stereotype. Intuitively, because the DM neglects non-representative types, he does not fully account the current observation may be due to sampling variability. As a consequence, his beliefs overreact when a type he does attend to is confirmed by the data. If criminal activity is part of a group’s stereotype, the DM over-reacts to seeing a criminal from that group and his judgments become even more biased against the group. If a growth company generates surprisingly positive earnings, investors further upgrade their belief that the stock is a good investment, because they neglect the possibility that an extreme observation may be due to noise.

At the same time, case ii) shows that the DM under-reacts (relative to a Bayesian) to information inconsistent with the stereotype. This is because insofar as the stereotype is unaffected, the probability of a non-stereotypical type is not upgraded, as the type remains neglected in the assessment of the group. This is the main departure of our model from the kernel of truth logic: the agent discards some news because these are not strong enough to overturn his prior belief.

Upon observing a highly successful member of a group stereotyped as low socioeconomic status, DMs code the occurrence as an “anomaly” and continue to believe that the group at

large should be viewed through the lens of the negative stereotype. People can espouse racist views and yet be friendly with individual members of the group they disregard (Schneider 2004). However, as shown in Proposition 8, non-stereotypical information is often ineffective at changing beliefs even if it swamps the few instances underlying the stereotype.³⁷

Proposition 8 implies that the DM exhibits a type of confirmation bias (Lord, Ross and Lepper 1979, Nickerson, 1998). Faced with two observations of different types from group G (formally, $n_{t,G} = n_{t',G} = 1$ and $n = 2$), such that t belongs to the stereotype of G but t' does not, the DM over-reacts to information consistent with the stereotype and ignores information inconsistent with it.³⁸ In this way, our approach provides a unified mechanism that gives rise to both base-rate neglect and confirmation bias: base-rate neglect arises when representative types are unlikely, while confirmation bias arises when new information does not change representativeness and allows stereotypes to persist. In the context of representativeness-based predictions, these biases are two sides of the same coin. The approach can also unify several other biases, such as overconfidence but also – under appropriate extensions – polarisation effects.³⁹

³⁷Propositions 7 and 8 formalise features of psychological models of stereotype change, such as the conversion and sub-typing models (Rothbart 1981, Weber and Crocker 1983), in which disconfirming evidence is treated as “exceptions to the rule” up to the point when it becomes sufficiently pressing, and engenders a stereotype change. In the experimental psychology literature on stereotype change, types are usually unordered and multidimensional (e.g., ethnicities are defined in terms of socio-economic status, proneness to violence, musical tastes, etc), see Hewstone et al (2000) and references therein. This literature finds that stereotypes fail to change when some extreme disconfirming evidence is observed (i.e. group elements that violate the stereotype along every dimension). In contrast, when sufficiently many group members are observed that violate the stereotype only along a given dimension, but are otherwise representative, those observations may become representative of the group as a whole and change the stereotype. Consistent with Proposition 7, stereotype change in these experiments seems driven by changes in (relative) frequency, and not by the shock of observing extreme exceptions.

³⁸Lord, Ross and Lepper (1979) suggest that confirmation bias arises in response to information that provides ambiguous support to an underlying hypothesis. Here, the information provided received presents ambiguous support for the stereotype.

³⁹Polarization arises as a consequence of confirmation bias when DMs have heterogeneous priors. Proposition 7 then implies that a given set of observations can lead different DMs with different stereotypes to each reinforce their own stereotype, and thus update in opposite directions.

G Experiments

G.1 Analysis of All Unordered Types Experiments

We conducted four experiments on unordered types. The final experiment, using cartoon characters in T-Shirts, were reported in the main text. Here we discuss the other experiments and their results.

G.1.1 Unordered Types Experiment 1: (Lots of) Triangles, Squares, and Circles

The first unordered types experiment used groups of 50 shapes each. The groups were characterized by color (red shapes or blue shapes) and the types were shapes (triangles, squares, and circles). In both conditions, the blue group contained 22 squares, 24 circles, and 4 triangles. In the Control condition, this blue group was presented next to a similar red group that contained 26 squares, 20 circles, and 4 triangles. Note that in the Control condition, within each group, the most representative type and the modal type coincide: among the blue shapes, circles are both most representative and modal, and among the red shapes, squares are both most representative and modal. In the Representativeness (Rep.) condition, we drive a wedge between modality and representativeness by changing the distribution of red shapes presented next to the blue group. In the Rep. condition, the red group contains 21 squares, 16 circles, and 13 triangles. While circles are still most representative and modal among the blue group, in the red group the modal shape is a square while the most representative shape is a triangle. Our prediction is that participants will be more likely to guess that the triangle is modal among the red shapes in the Rep. condition than in the Control condition. The images as they appeared to participants are reproduced in Figure 7.

This design is not as clean as the T-shirts design presented in the paper. Most importantly, if we do see an increase in the fraction of participants that believe triangles are modal in the Rep. condition, we cannot rule out that this is simply driven by the fact that there are more red triangles in the Rep. condition than in the Control condition. We present the results below with that caveat.

This shapes experiment was conducted on MTurk in November 2014 with 217 partici-

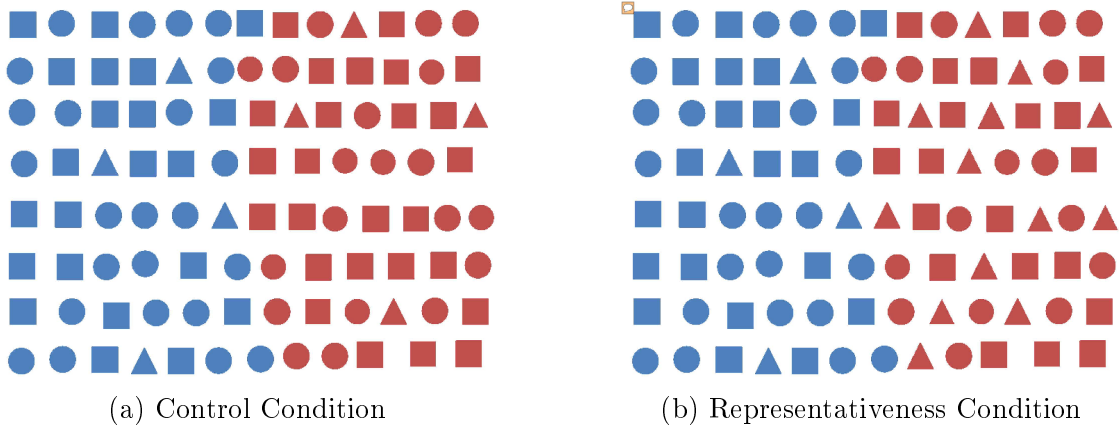


Figure 7: 50 Shapes Experiment

pants.⁴⁰ Participants viewed the shapes for 15 seconds and then completed 10 simple addition problems (computing sums of two-digit numbers) before answering a series of questions about the shapes they saw. They were asked to guess what the most common shape among each group was and to estimate the frequency of each shape in each group. They received \$0.30 for completing the HIT and an additional \$1 if they answered one of the randomly-selected questions about the shapes correctly.

We find that 7% of participants in the Control treatment and 13% of participants in the Rep. treatment believe that the triangle is the modal red shape. The direction matches our prediction but the effect is not significant at conventional levels ($p=0.17$).

After conducting this experiment, we altered the design to eliminate the potential confound. In all designs going forward, we hold fixed the number of objects of the type of interest across the Control and Rep. treatment and simply alter the comparison group to change whether or not the type is diagnostic.

G.1.2 Unordered Types Experiment 2: (Fewer) Triangles, Squares, and Circles

The next iteration improved on the original shapes design in a few important ways. First, we cut down the number of shapes, reducing the groups from 50 shapes each to 25 shapes each. Second, we changed the distributions such that the number of red triangles was held

⁴⁰This count excludes 3 participants who self-identified as color blind. Neither the point estimates or p-values reported below are changed if those participants are included in the analysis. The HIT was posted once for 200 participants and we had 220 complete the experiment on Qualtrics via the link (some fail to submit the payment code to MTurk for payment, allowing us to overshoot our target.)

constant across condition, but the number of blue triangles varied to change how diagnostic the triangles were for the red group. In both conditions, the red group contained 6 squares, 10 circles, and 9 triangles. In the Control condition, this group was presented next to a blue group that contained 9 squares, 8 circles, and 8 triangles. In the Rep. condition, the red group was presented next to a different blue group that contained 11 squares, 12 circles, and 2 triangles. Thus, while the number of red triangles is the same across conditions, triangles are much more representative of the red group in the Rep. condition than in the Control condition. We predict this shift in representativeness of the red triangles will lead to an increase in the proportion of participants who guess that triangles are modal in the red group and an increase in the estimated frequency of red triangles.

We ran this experiment both on MTurk and at the Stanford Experimental Economics Laboratory in January 2015. The MTurk protocol was very similar to Experiment 1, the previous shapes experiment. Participants viewed the objects for 15 seconds, answered 10 simple addition questions, then answered a series of questions about the shapes. Participants were paid \$0.30 for completing the HIT and an additional \$1 if they answered a randomly-selected question about the shapes correctly. We collected data from 100 participants, 50 in each condition.⁴¹

In the Control condition, 18% of participants believed triangles were modal in the red group; in the Rep. condition, this grows to 24% ($p=0.46$ from two-tailed test of proportions).⁴² Participants in the Rep. condition estimate that there are 9.98 red triangles on average, while participants in the Control condition estimate that there are 9.39 red triangles on average (two-tailed t-test, $p=0.65$). But, this difference is largely driven by one participant who provided an unusually large estimate of red triangles in the Rep. condition (50). If we exclude this participant, the data on estimated frequencies is not directionally consistent with our hypothesis, with the average estimate of red triangles being 9.39 in the Control condition being and 9.16 in the Rep. Condition (two-tailed t-test, $p=0.82$).

The protocol in the Stanford laboratory was more complicated, with several potentially

⁴¹This count excludes 1 participant who self-identified as color blind. Including this participant does not impact the results presented below. We posted the HIT once for 100 participants.

⁴²Using a probit regression that controls for demographics (gender and year of birth) also estimates approximately a 6 percentage point increase in the fraction of participants that believe the triangles are modal in the red group.

important changes. First, instead of arranging the shapes on a page for participants, we provided participants with an envelope that contained cutouts of each of the 50 total shapes for their condition. Participants were given 1-minute to open the envelope and view the contents. Second, in the laboratory, we had participants complete both an ordered and an unordered types experiment, back-to-back, in a randomly-assigned order. Third, after viewing the objects in the envelope and completing the math problems, participants were asked to describe their envelope, in writing, to another participant in the lab. This was incentivized as “advice”. Take a participant who had been given an envelope labeled “A” (i.e. was assigned to the Control condition). We told this participant that later in the experiment, we were going to ask another participant in the lab, who had been given a different envelope, a question about envelope “A”. This participant would receive the advice, but not the envelope. If the participant answered the question about envelope “A” correctly, both the advice giver and the other participant would receive additional payment. Thus, participants were incentivized to write down information about the shapes in their envelope that would be accurate and useful. Thus, we likely encouraged some careful reflection on their envelope before the participant had answered any of our other questions of interest about the shapes. We ran four laboratory sessions, with 66 total participants.⁴³

The Stanford laboratory results do not support our hypotheses. In the Control condition, 33% of participants believe triangles are the modal red shape; in the Rep. condition, 27% of participants believe triangles are the modal red shape ($p=0.59$ from two-tailed test of proportions). This result does not depend on whether participants completed this unordered types experiment first or second. Participants also estimate -0.61 fewer red triangles in the Rep. condition than in the Control condition. This difference goes in the opposite direction of our prediction, though it is not significant.

The results for this design are the weakest among our unordered types experiments. While we do not have conclusive evidence on what drives these effects, we do have a hypothesis that seems consistent with the data. It may be the case that participant judgments were swayed by the total number of each type, pooled across groups. Consider the triangle questions.

⁴³Our ex ante plan was to run four sessions, though we had thought this would yield closer to 100 participants. After four sessions, we stopped and attempted to improve the design as described below.

We expect that the 9:2 red triangle to blue triangle ratio in the Rep. condition, relative to the 9:8 red triangle to blue triangle ratio in the Control condition, will lead participants to estimate a larger share of red triangles. But, it is also true that they see 11 total triangles in the Rep. condition, but 17 total triangles in the Control condition. In the laboratory experiment, unlike on MTurk, the shapes are not arranged by group for participants; they are loose in an envelope. If the distinction between the groups is not natural at the moment when they are forming their impressions of the envelope they saw, the fact that there are fewer total triangles may carry more weight than the representativeness of the triangle within each group. This force may push them toward estimating that there were fewer red triangles in the Rep. condition.

After these results, we sought to improve the experiment. In particular, we moved to simpler distributions, where only two types of objects appeared within a given group. This amplified the extent of diagnosticity, as certain types now appear in only one of the two groups. We also shifted from using shapes to more familiar objects, thinking that this might make “groups” a more natural concept. We also switched back to displaying the objects in a fixed arrangement for participants, so we could arrange the objects into obvious groups.

G.1.3 Unordered Types Experiment 3: Cars, Trucks, and SUVs

Next, we ran a version of the experiment that used groups of vehicles. The groups were defined by color, with a group of blue vehicles and a group of green vehicles. The types were defined by type of vehicle: pick-up truck, sedan, or SUV. Each group had 20 vehicles. The distributions were similar to the T-shirt design. The green group of vehicles consisted of 9 SUVs and 11 sedans. In the Control condition, this group was displayed next to a group of blue vehicles with the same distribution, 9 SUVs and 11 sedans. Thus, in the Control condition, there is no vehicle type that is diagnostic of a group. In the Rep. condition, the green group was displayed next to a blue group with 9 trucks and 11 sedans. As with the T-shirts design, this creates a tension between the modal type and the diagnostic type in each group. In the green group, the sedan is modal but the SUV is diagnostic; in the blue group, the sedan is modal but the truck is diagnostic. Thus, we predict that participants in the Rep. condition will be more likely to guess that the “9 vehicle” type is modal for a

group, because the 9-vehicle type is diagnostic in this condition. The images, exactly as they appeared to participants, are reproduced in Figure 8.

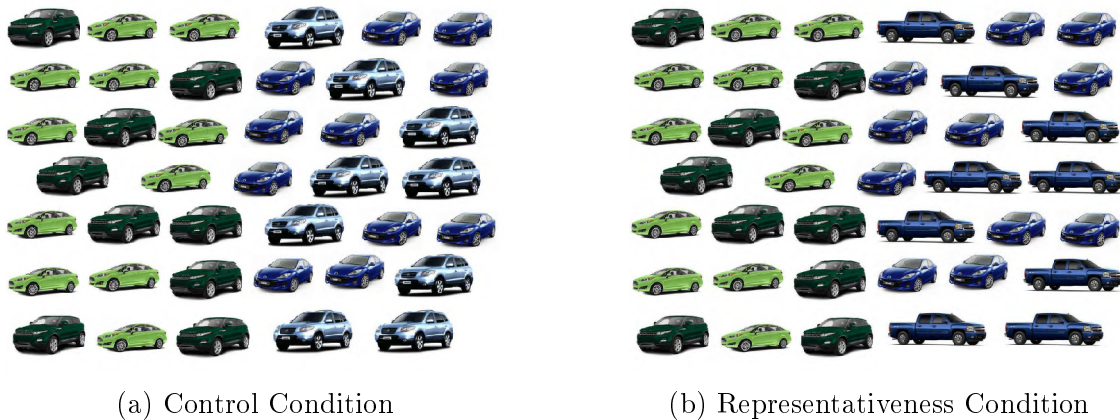


Figure 8: Vehicles Experiment

We conducted this experiment with 57 participants on MTurk in January 2015.⁴⁴ The protocol was very similar to the T-shirts experiment reported in the main text. Participants were given 15 seconds to review the objects, seeing the green group next to a randomly-chosen comparison group, either the Control blue group or the Rep. blue group. Then, participants were asked what the most common type of vehicle was for each group and were asked to estimate the frequency of different types of vehicles for each group. Participants received \$0.30 for completing the HIT and an additional \$2 in incentive pay if they answered a randomly-selected question correctly.

Our results support our hypothesis. In the Control condition, when the 9-vehicle type is not diagnostic, participants guess that the 9-vehicle type is modal in 22% of cases. In the Rep. condition, when the 9-vehicle type is diagnostic of each group, participants guess that this 9-vehicle type is modal in 40% of cases (significantly different than the Control condition using a two-tailed test of proportions with $p = 0.042$). Because we have two observations per individual (her guess of the most common blue vehicle type and her guess of the most common green vehicle type), it is useful to run a probit regression that allows us to cluster observations at the individual level. When we predict the probability of guessing the 9-

⁴⁴The HIT was posted once, for 150 participants, to be randomized in equal proportions into this experiment and the ice cream ordered types experiment. We collected data from 76 participants for this experiment, but 19 had participated in a previous version of the experiment, leaving us with 57 participants. The results below are directionally stronger if we include those repeat participants.

vehicle type is modal from a participant’s randomly-assigned treatment, her demographic information (gender and year of birth), and a dummy for whether the guess was for the blue or green vehicles, we estimate that participants in the Rep. condition are 17.4 percentage points more likely to guess the 9-vehicle type is modal ($p=0.09$).

We can also look at estimated frequencies of different types across condition. In this experiment, we only asked participants to estimate the number of green sedans and SUVs and blue sedans and pickups (the types that appeared in the Rep. conditions). Thus, because we are missing estimates of the blue SUVs, we cannot do quite the same analysis presented for the T-shirts design, where we compared the estimate of the modal type and the 9-vehicle type for each group across conditions. But, we can do this analysis for the green group, asking how the estimated difference in number of sedans and SUVs varies across conditions. We predict that participants will estimate a greater gap between green sedans and SUVs in the Control than in the Rep. condition. Using an OLS regression, we predict the estimated difference between the number of green sedans and green SUVs from a participant’s randomly-assigned condition and her demographic information. We find that the effect is small but directionally supportive of our hypothesis, with participants in the Rep. condition estimating the difference in sedans and SUVs to be about 0.5 counts smaller than participants in the Control condition ($p=0.68$).

We moved from using the vehicles to using the cartoon characters wearing T-shirts in an attempt to simplify the objects. The pictures of vehicles are highly detailed, providing many features that could capture participants’ attention during their brief 15-second viewing. Furthermore, recognizing the same type across group was not straightforward – i.e. the green sedan and blue sedan have many differences in addition to color. We wanted to move to a format where fewer features would vary, and where recognizing the same type across group would be simpler. This led us to the T-shirts design.

G.1.4 Unordered Types Experiment 4: T-Shirts

The T-Shirts design was reported in the main text. We ran this experiment in the laboratory and on MTurk. On MTurk, participants received \$0.30 for completing the experiment and an additional \$1 if they answered the randomly-selected question correctly. Data was collected

in February 2015. The laboratory sessions were conducted at the Ohio State Experimental Economics Laboratory in March 2015. Participants dropped into the lab for approximately five minutes, receiving a \$5 show-up fee and up to \$5 more in incentive pay. In the lab, we added two questions on risk preferences between the viewing of the objects and the questions about the T-shirt people in order to better obscure our focus.

We had 301 total participants, 196 in the laboratory and 105 on MTurk.⁴⁵ We have two observations for each individual: her guess of the most common color shirt among the girls and her guess of the most common color shirt among the boys. Our main hypothesis is confirmed in the pooled data (including guesses about both girls and boys): participants in the Control condition believe the 12-shirt color is modal in 35% of cases, while this mistake is made in 46% of cases in the Rep. condition ($p=0.01$ from two-tailed test of proportions). Using a probit regression that clusters observations at the individual level, we estimate that when the 12-shirt color is diagnostic of a group, a participant is 10.5 percentage points more likely to believe it is the modal color ($p=0.01$). This effect is significant when we restrict attention to the sample from the laboratory (14.4 percentage points, $p=0.007$) and directional in the smaller MTurk sample (7.6 percentage points, $p = 0.26$).

We also analyze the difference in estimated counts of the modal color shirt and the counts of the 12-shirt color shirt the participant saw (we subtract estimated counts of the 12-shirt color from estimated counts of the modal color for each participant for each group). We find that, on average, participants in the Control condition estimate having seen 0.54 more modal color shirts than 12-shirt color shirts, while participants in the Rep. condition estimate having seen 0.72 fewer modal color shirts than 12-shirt color shirts (this across treatment difference is significant with $p = 0.013$ using a two-tailed Fisher Pitman permutation test). Using an OLS regression, we find that when the 12-shirt color is representative, participants estimate the difference in counts between the true modal color and the 12-shirt color to be 1.39 counts smaller ($p=0.006$). The results are similar and significant within either subsample, lab or MTurk.

⁴⁵We recruited 150 participants for the MTurk experiment, but 45 who completed our HIT had already completed a previous version of the experiment and are excluded from our analysis. The target for the laboratory sample was 200 participants over three days of drop-in sessions. We had 202 participate, but we exclude 6 laboratory participants who self-reported color blindness. The results are very similar if all of these participants (both repeat participants for MTurk and color blind) are included.

G.1.5 Summary of Unordered Types Experiments

Table 3 summarizes the results from the four unordered types designs. For each experiment, we run a probit regression predicting the probability that the participant believed a less common type was the modal type from whether or not the type was representative. For the vehicle and T-shirts experiments, we have two observations per individual and we cluster the standard errors at the individual level. For the shapes experiments, we have one observation per individual. We report the marginal effect of assignment to the Rep. condition (where the less common type was representative) on the probability of guessing that the less common type was modal. The last row reports the same coefficient, but from a probit regression that uses all of the data from the unordered types experiments. We include a dummy for each particular experiment and cluster observations at the individual level. We find a directional effect consistent with our hypothesis in five of the six samples – all but the Stanford laboratory sample for Experiment 2 (the 25 Shapes design). When we pool all data, we estimate that a participant is 9.3 percentage points more likely to believe the less common type is modal when it is representative than when it is not ($p=0.002$). If we include, in addition, all color blind participants, this estimate is 9.0 percentage points ($p=0.002$); and, if we include all observations, including all observations from participants who have participated in previous versions of the experiments, this estimate is 8.3 percentage points ($p=0.003$).

We perform a similar analysis using the data on estimated frequencies. The ideal analysis would look at how the magnitude of the difference in the estimated frequency of the modal type less the estimated frequency of the less common type changes across condition. We can do this calculation in Experiments 2 - 4. In Experiment 1, the true frequency of the less common type varies across condition, so this analysis is not useful. In Experiments 2 - 4, the true frequency of both the modal type and the less common type are held constant across treatment. Therefore, we can explore how this difference varies based upon whether the less common type is representative. The prediction is that the difference in estimated frequencies should decrease when the less common type is representative, as participants in the representativeness condition will estimate fewer counts of the modal type (as it is now

less representative for the group) and more counts of the less common type (as it is now more representative for the group).⁴⁶ When we pool the data from Experiments 2 - 4, we estimate that the difference in estimates of the modal type and the less common type decrease by approximately 1.06 counts in the representativeness condition ($p=0.010$). If we include, in addition, color blind participants, this estimate is a decrease of 1.09 counts ($p=0.008$), and if we include all observations, including those from participants who have participated in multiple versions of the experiment, the estimate is a decrease of 1.19 counts ($p=0.002$).

G.2 Analysis of All Ordered Types Experiments

We conducted two experiments on ordered types. The final version, using ice cream cones, was reported in the main text. Here, we report the other experiment and discuss the complete set of results.

G.2.1 Ordered Types Experiment 1: Rectangles

Our first design for the ordered types experiment used groups of rectangles of varying heights. We created a group of blue rectangles, each of which were 1-unit wide and 1, 2, 3, 4, or 5 units tall. In the Control condition, this group was presented next to a group of red rectangles of the same width with a very similar distribution over heights. In the Rep. condition, the blue group was presented next to a red group of rectangles with the same width, but with a distribution over heights that created a representative tall type for the red group. Table 4 displays the distribution, and Figure 9 presents the images, exactly as they appeared to participants.

In the Control condition, no type is very representative of either group, and the small difference in the distributions occurs at types close to the mean. In the Rep. condition, on the other hand, we create a highly representative type for the red group, as there are five

⁴⁶Note that in Experiment 2, we did not ask participants for their estimates of counts of the modal type. Therefore, we simply analyze the change in the difference 0 minus the estimated counts of the less common type for that experiment. We have one observation for each individual (0 - estimate of red triangles). For Experiment 3, we also have one observation per individual (estimate of green sedans - estimate of green SUVs). For Experiment 4, we have two observations per individual (estimate of modal color - estimate of 12-shirt color for both boys and girls). We cluster at the individual level, giving us 824 observations for 523 individuals.

Table 3: Summary of All Unordered Types Experiments

| Experiment | # of Subjs. | Percentage Point Increase in Prob. Guess Less Common Type is Modal when it is Representative | p-value | Change in Diff. in Estimated Frequencies $\Delta(\text{Modal} - \text{Less Common})$ | p-value |
|-----------------------|-------------|--|--------------|--|--------------|
| 1: 50 Shapes on MTurk | 217 | 5.6 pp | 0.17 | N/A | |
| 2: 25 Shapes Pooled | 166 | 1.4 pp | 0.84 | -0.12 | 0.88 |
| MTurk Only | 100 | 5.6 pp | 0.49 | -0.50 | 0.70 |
| Lab Only | 66 | -13.8 pp | 0.28 | 0.61 | 0.34 |
| 3: Vehicles on MTurk | 57 | 17.4 pp | 0.09 | -0.45 | 0.68 |
| 4: T-Shirts Pooled | 301 | 10.5 pp | 0.014 | -1.39 | 0.006 |
| MTurk Only | 105 | 7.6 pp | 0.25 | -2.19 | 0.012 |
| Lab Only | 196 | 14.4 pp | 0.007 | -1.16 | 0.056 |
| Pooled | 741 | 9.3 pp | 0.002 | -1.07 | 0.010 |

Notes: Std. errors are clustered at the individual level. We report the marginal effect of the coefficient on treatment from a probit regression predicting the probability of the error. Each specification includes all demographic variables collected for that experiment. The pooled specification includes only treatment and gender, as this is the only demographic variable that was collected across all experiments.

Table 4: Distributions for Ordered Types Experiment 1

| Height in Units (Types) | Counts for Blue Group | Counts for Control Red Group | Counts for Rep. Red Group |
|-------------------------|-----------------------|------------------------------|---------------------------|
| 1 | 3 | 3 | 4 |
| 2 | 8 | 9 | 11 |
| 3 | 24 | 23 | 20 |
| 4 | 14 | 14 | 10 |
| 5 | 1 | 1 | 5 |
| Total Counts | 50 | 50 | 50 |
| Mean Height | 3.04 | 3.02 | 3.02 |

5-unit tall rectangles in the Rep. red group and only one 5-unit tall rectangle in the blue group. Importantly, across both conditions, the means of the two groups are held constant, with the blue group always having a mean height of 3.04 units and the red group having a mean height of 3.02 units. The prediction is that participants will be more likely to guess that the red rectangles are taller on average in the Rep. Condition than in the Control condition, because of the representative tall type among the Rep. red group.

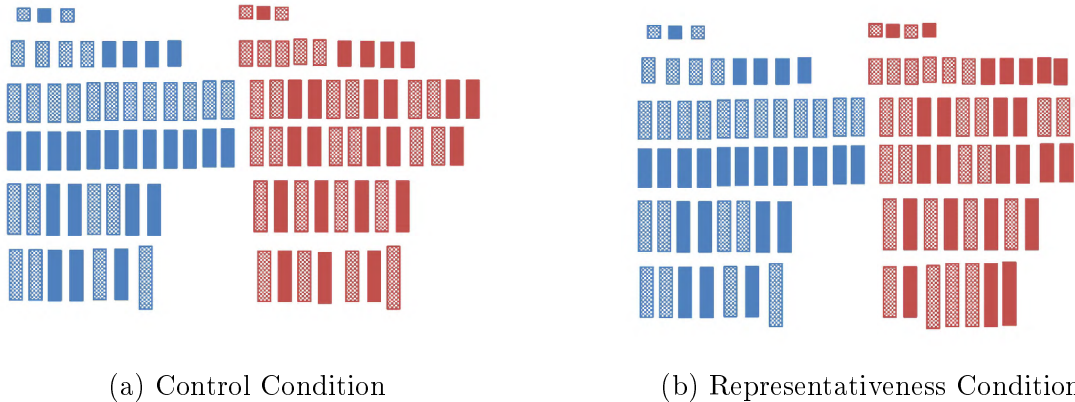


Figure 9: Rectangles Experiment

We chose to arrange the rectangles by height for participants so that it might be easier to digest and make sense of the groups in a short period of time. We also varied the fill of the rectangles, with half of each group’s rectangles having a solid fill and half displaying a checkered fill.⁴⁷ Our fear was that if only the heights varied by shape, participants might anticipate that we were particularly interested in their impressions of the heights of the rectangles. So, we chose to vary the fill as well to create another plausible dimension of interest.

The first experiment using this rectangles design was conducted on MTurk in November 2014 with 113 participants.⁴⁸ Participants were randomly-assigned to view either the Control rectangles or the Rep. rectangles for 15 seconds. Then, they completed simple addition problems for approximately 3 minutes, computing sums of two-digit numbers. Finally, par-

⁴⁷The fill was performed such that approximately half of each type within each group received each fill. That way, the representativeness patterns we sought to induce in the distributions were preserved within each fill.

⁴⁸The HIT was posted once for 100 participants, and 114 completed the experiment via the link to Qualtrics. We exclude one participant who self-identified as color blind. The point estimates and p-values reported below are unchanged if this participant is included.

ticipants were asked questions about the shapes they saw, including which color rectangles were taller on average, which group of rectangles they would prefer to choose from if they were going to earn \$0.50 per unit height of a randomly-drawn rectangle, and the average height of each group of rectangles. We also asked about the fill of the rectangles they saw, so not all questions would focus on height. Participants received \$0.30 and up to an additional dollar in incentive pay based upon their answers to the questions about the shapes.

The results are consistent with our hypotheses. In the Control condition, 40% of participants believed the red group was taller on average, while in the Rep. condition, 60% of participants believed the red group was taller on average ($p=0.03$ from two-tailed test of proportions). When we look at which group of shapes participants preferred to bet on, the results are weaker but still directionally supportive: 45% of participants in the Control and 58% of participants in the Rep. group prefer to choose from the red shapes when they will be paid based upon the height of a randomly-drawn rectangle ($p=0.16$ from two-tailed test of proportions).

There is no difference in estimated average height of the red shapes across condition (3.28 in the Control versus 3.29 in the Rep. condition, $p=0.95$). If we look at the estimated average height of the blue shapes – recall that the blue shapes are identical across condition – we see that participants in the Rep. condition believe they are slightly smaller on average, though this difference is not significant (3.30 in the Control versus 3.22 in the Rep. condition, $p = 0.59$).

We took this design into the laboratory in January 2015 at the Stanford Experimental Economics Laboratory. There were a few potentially important changes to the protocol in the laboratory. For one, we had participants complete both an ordered and an unordered types experiment, back-to-back, in a randomized order. Note that this is the same sample for whom we reported results for the unordered types Experiment 2 above. Instead of participants viewing the objects on a computer screen, we passed out envelopes that contained a printed handout of either the Control or the Rep. shapes. After viewing the handout in the envelope and completing the math problems, participants were asked to describe the handout they had seen, in writing, to another participant in the lab. This was incentivized as “advice”, implemented as described in the previous Stanford laboratory description for unordered types

Experiment 2. We ran four laboratory sessions, with 66 total participants.

The results from the laboratory were inconsistent with our hypotheses. We find that 46% of participants in the Control condition believed the red shapes were taller on average, while only 36% of participants in the Rep. condition made this error ($p=0.45$ from two-tailed test of proportions). When we look at the choices about which group participants preferred to bet on, the results are even more striking. Nearly 67% of participants in the Control condition prefer to choose from the red shapes, while only 27% of participants in the Rep. condition prefer to choose from the red shapes ($p=0.001$ from two-tailed test of proportions). Looking at the data on estimated average heights across condition, there are no significant differences. Directionally, participants estimate both the blue and the red shapes to be taller on average in the Rep. condition than in the Control condition.

There are a few issues with the rectangles design that we sought to address in later experiments. First, it may have been tricky for participants to recognize and process heights of rectangles. We tried to describe the types in terms of “units” of height, but this likely felt a bit confusing to participants. Therefore, we wanted to move to an ordered space that had more obviously distinct types. That is why we shifted to using “scoops” of ice cream, where the difference between 1, 2, 3, 4, or 5 “units” would be more easily recognizable and familiar. Second, there may have been too many shapes on the page for participants to make sense of in a 15-second viewing period. Looking at the advice participants wrote in the laboratory sessions is very informative. Many participants accurately recalled and described the first row of rectangles (featuring three 1-unit tall blue rectangles and four 1-unit tall red rectangles in the Rep. condition, and three 1-unit red and blue rectangles in the Control condition), but no advice sheet even attempted to describe the final row. It may be that with only 15 seconds, participants only have time to focus on part of the page, and the top of the page may be a likely place to start. This type of behavior would hurt us substantially: if participants are mostly focused on the top of the page, they will miss out on the representative tall types we generated. Even worse, in the first row, there are more short red shapes than short blue shapes in the Rep. condition but not in the Control condition. This could lead to participants thinking, contrary to our prediction, that the red group is shorter on average in the Rep. condition. If the first row or two is what participants

mainly recall, it could also explain why so many participants prefer to bet on blue in the Rep. condition, as they remember there were more of the worst possible payoff shapes among the red group. We decided to cut down the number of objects in order to give participants a better chance to view the group as a whole during a short window. And, perhaps more importantly, we altered the distributions so that the group with the representative tall type would not also have comparatively more of the shortest possible type.

G.2.2 Ordered Types Experiment 2: Ice Cream

After running the rectangles experiments, we sought to simplify the protocol as much as possible. We did this by reducing the number of objects, but also by eliminating the math problems from between the viewing of the objects and the answering of our questions of interest. This led to the ice cream cone design, illustrated in Figure 10.

Groups are sets of 24 ice cream cones: group membership is defined by ice cream flavor (chocolate vs strawberry), and types are the number of ice cream scoops, ranging from 1 to 5. In the Control condition, Fig.10a, distributions are very similar, with most cones having intermediate numbers (2 or 3) of scoops. Here, no type is particularly representative of either group. In the representativeness condition, Fig.10b, the same chocolate cones are presented next to a different group of strawberry cones. In the Representativeness condition, strawberry cones have the same average number of scoops as do the Control condition strawberry cones, but, importantly, they do not contain any 5-scoop cones. This makes the right tail, 5-scoop cones very representative for the chocolate group. Similarly, in this condition only the strawberry group has a cone with 1 scoop, making the left tail very representative for that group.

We ran this ice cream cone design on MTurk in January 2015 with 65 participants.⁴⁹ When asked which flavor had more scoops on average, 34% of the Control condition guesses chocolate and 67% of the Rep. condition guesses chocolate ($p=0.009$). We ask participants a related question using choices over lotteries. They are told that we are going to randomly

⁴⁹We posted the HIT once for 150 participants, with participants randomized in equal proportions into either this experiment or the vehicles experiment described above. Eighty-four MTurk participants completed this experiment, but 19 of these individuals had participated in a previous version of this experiment and therefore are excluded from this analysis. The results reported are unchanged if these participants are included.

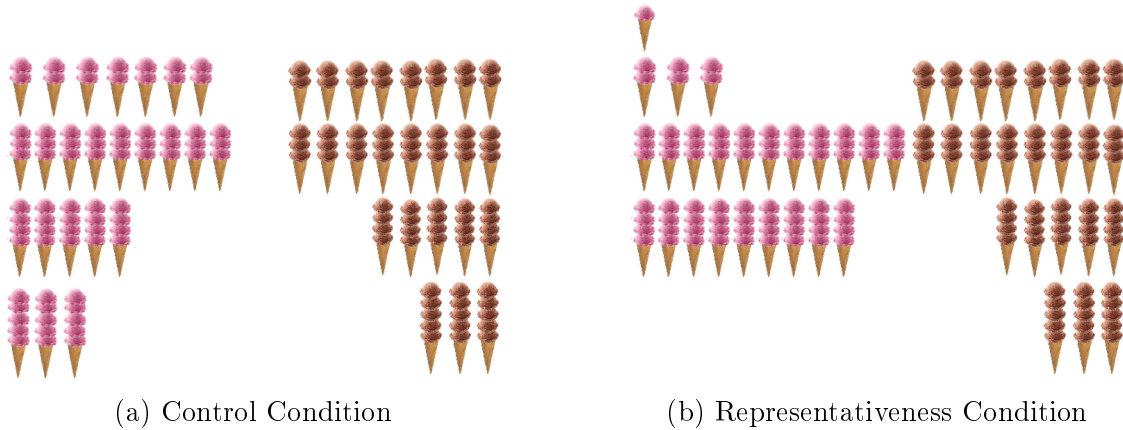


Figure 10: Ice cream cone Experiment

choose one of the ice cream cones they saw, with the participant earning \$0.50 for every scoop the randomly-chosen cone has. They are asked to choose which flavor we draw from. The proportion that chooses the chocolate lottery grows from 37% in the Control to 57% in the Rep. condition ($p=0.11$). Finally, we explore the participants' estimates of the average number of scoops on both the chocolate and strawberry cones. In the Control condition, participants believe the strawberry cones have on average 2.85 scoops and the chocolate cones have on average 2.82 scoops. In the Rep. condition, participants believe the strawberry cones have 2.82 scoops on average and the chocolate cones have 2.71 scoops on average. None of these differences, either across condition or flavor, are significant. Overall, the fraction of participants who provide greater estimates of the average number of chocolate scoops than the average number of strawberry scoops is larger in the Rep. conditions than in the Control conditions (60% versus 40%, $p = 0.11$ from two-tailed test of proportions).

After running this experiment on MTurk, we sought to bring this design into the laboratory. The first ice cream laboratory protocol used the same ice cream cone images with participants directed to answer our questions of interest immediately following the viewing of the objects. In this way, it was likely quite clear to participants what our goal as researchers was: to test their recall of the images they saw. While participants on MTurk are often asked simple attention checks or to report basic objective information (who is this a picture of, transcribe this audio clip, answer this survey question), participants in the laboratory are likely less familiar with this type of design. It is possible that they were skeptical or wary of

being tricked – i.e. why am I being asked what seems to be an obvious question?

We had 56 participants complete this experiment at the Ohio State Experimental Economics Laboratory in March 2015 before we stopped to evaluate what was going on. The results from this experiment look similar to the laboratory data from Stanford. In the Control condition 45% of participants believe the chocolate cones had more scoops on average, while only 36% of participants in the Rep. condition make this error. The effect goes in the opposite direction of our prediction, but is not significant ($p = 0.48$ from two-tailed test of proportions). Similarly, the proportion of people who prefer to bet on the chocolate cones falls from 48% in the Control to 32% in the Rep. condition ($p = 0.21$). There is no difference in estimated average number of scoops of either flavor across condition.

Having seen these results, we brainstormed why participants in the Rep. condition in the lab would be less likely to believe the chocolate cones are taller on average. A directional effect opposite the predicted direction suggests something else at work – something that is not at work on MTurk, where both the rectangle and the ice cream design produced results that support our hypotheses. We conjecture that participants in the laboratory are more skeptical of being tricked, perhaps because they are not usually asked something so simple in a typical economics experiment. It may also be that our ice cream design was “too good” – that is, from a quick look at the objects, the chocolate cones quite strikingly appear to have more scoops on average, that participants are worried that this is actually a trick question. We do not have direct data on this issue, but we did change the design in an attempt to address this problem head on.

We used the same distributions of ice cream cones, but added a new, small section on risk preferences between the viewing of the objects and the questions about the objects. This creates a plausible alternative research question – we could be interested in how viewing a particular arrangement of ice cream cones impacts a participant’s risk preferences. We paid participants for this risk preference section, and we framed the questions about the ice cream cones as more of an attention check than an item of interest. We also added a question to the very end of the experiment asking participants what they believed the experiment was trying to test. In this new design, no participant correctly identified our focus on number of scoops. Our interpretation of the data is that the introduction of this “decoy” encourages

less skepticism on the part of participants, and perhaps helps us more successfully elicit their quick, gut reactions to the objects, much the way we were able to do on MTurk. This is speculative, but it does seem consistent with the data we have collected.

We had 101 new participants from the Ohio State Experimental Economics Laboratory complete the updated ice cream protocol.⁵⁰ When asked which flavor had more scoops on average, 51% guess chocolate in the Control and 56% guess chocolate in the Rep. condition ($p=0.61$). There are no significant differences in estimates of average scoops across flavor or condition. The Rep. treatment produces an insignificant decrease in the proportion that prefer the chocolate lottery (45% to 38%, $p=0.47$).

A natural question to ask is why results for the choices over lotteries would be weaker than the results for which flavor had more scoops on average (or which shapes were taller in the rectangles experiment). While an individual who believes that chocolate cones have more scoops on average should believe there is also a greater expected value from the chocolate cone lottery, it does not guarantee that the chocolate cone lottery is the expected utility maximizing choice: risk preferences may also play a role. Therefore, indicating that chocolate cones have more scoops on average does not guarantee that a reasonable participant will also choose the chocolate cone lottery. To shed light on this issue, we asked a different set of participants from the same laboratory population about their hypothetical preferences over these lotteries. We presented the three lotteries (chocolate cones, Control strawberry cones, and Rep. strawberry cones) side-by-side, described as abstract gambles (there was no mention of ice cream and no visual representation of the lotteries). They were then asked to rank the attractiveness of these gambles from most to least attractive. In a sample of 196 participants, 22% prefer the chocolate cones lottery to the lottery induced by the Rep. strawberry, while 39% prefer that same chocolate cones lottery to the lottery induced by the Control strawberry cones. This suggests that risk preferences were likely working against us finding an effect in support of our hypothesis, as this data would predict a 17 percentage point decrease in the proportion choosing chocolate under the Rep. condition. In light of this baseline, the fact that we see only a 10 percentage point decrease in the lab and a 20 percentage point *increase* on MTurk suggests that the presence of diagnostic types is shifting

⁵⁰Our ex ante target was 100 participants over two days of drop-in sessions.

choices in line with our hypothesis.

G.2.3 Summary of Ordered Types Experiments

Table 5 summarizes the results from the two ordered types experiments. For each experiment, we run a probit regression predicting the probability that the participant guessed the shorter group was taller on average from her treatment assignment. We report the marginal effect of assignment to the Rep. condition (where the tallest possible type is the most representative type in the shorter group) on the probability of guessing that the shorter group is taller on average. The last row reports the same coefficient, but from a probit regression that uses all of the data from the ordered types experiments. We include a dummy for each particular sample and cluster observations at the individual level. We find a directional effect consistent with our hypothesis in three of the five samples. When we pool all data, we estimate that a participant is 9.3 percentage points more likely to believe that the shorter group is taller on average when it has a representative tall tail ($p=0.062$). If we include, in addition, all color blind participants, this estimate is 9.6 percentage points ($p=0.055$); and, if we include all observations, including participants who have participated in multiple versions of the experiment, this estimate is 8.4 percentage points ($p=0.083$). Note that for ordered types experiments, there is a significant difference between the laboratory studies and the MTurk studies. Using only the MTurk example, we estimate that a participant is *25 pp* more likely to guess that the shorter group is taller when it has a representative tall tail ($p=0.001$); the estimate for the laboratory sample is directionally negative, -3.1 pp ($p=0.65$). This difference in treatment effect across platform is significant ($p=0.006$).

We also provide the estimated treatment effect on the probability of choosing the shorter group to bet on for each experiment and the pooled estimate. There is no support for this prediction in the data (see discussion on confound of risk preferences above).

Table 5: Summary of All Ordered Types Experiments

| Experiment | # of Subjs. | Percentage Point Increase in Prob. Believed Shorter Group is Taller in Rep. Condition | | Percentage Point Increase in Betting on Shorter Group in Rep. Condition | |
|----------------------|-------------|---|--------------|---|-------------|
| | | | p-value | | p-value |
| 1: Rectangles Pooled | 179 | 7.3 pp | 0.32 | -6.8 pp | 0.37 |
| MTurk Only | 113 | 19.1 pp | 0.04 | 12.8 pp | 0.18 |
| Lab Only | 66 | -11.0 pp | 0.38 | -39.3 pp | 0.001 |
| 2: Ice Cream Pooled | 223 | 9.2 pp | 0.17 | -1.7 pp | 0.80 |
| Lab, No Decoy | 56 | -12.0 pp | 0.37 | -14.1 pp | 0.28 |
| MTurk | 65 | 30.7 pp | 0.01 | 18.8 pp | 0.13 |
| Lab, Decoy | 101 | 3.5 pp | 0.73 | -7.5 pp | 0.45 |
| Pooled | 402 | 9.3 pp | 0.062 | -3.8 pp | 0.45 |

Notes: Std. errors are clustered at the individual level. We report the marginal effect of the coefficient on treatment from a probit regression predicting the probability of the error. Each specification includes all demographic variables collected for that experiment. The pooled specification includes only treatment and gender, as this is the only demographic variable that was collected across all experiments.

H Empirical Analysis: Further Results

We repeat the analysis in the main text, implementing the regressions in Equations (6, 7). But, we add an additional control: the average likelihood of tail positions. This is the average frequency of the three types above the median for conservatives and the average frequency of the three types below the median for liberals. We again test the hypothesis that Θ is a significant predictor of $\mathbb{E}^{st}(t|cons)$ with a positive sign, and a predictor of $\mathbb{E}^{st}(t|lib)$ with a negative sign. Table 6 shows that, conditional on true mean and our measure of likelihood of tail positions, Θ predicts believed mean for each group G as predicted.

H.1 Additional Tables on Model Predictions

In this section, we present additional details on the model predictions. Table 7, Panel (a) summarizes the results for the GNH dataset, reporting not only MSPE, but also mean prediction error (MPE) and the fraction of observations for which the model underestimates

Table 6: Θ Predicts Beliefs Controlling for Likelihood

| | OLS Predicting Believed Mean of Group G | | | | | |
|---------------------|---|------------------|--------------------|--------------------|------------------|---------------------|
| | G = Conservatives | | | G = Liberals | | |
| | GNH | ANES | Pooled | GNH | ANES | Pooled |
| True Mean of G | 0.15 (0.36) | 0.65 (0.46) | -0.10 (0.43) | 0.92**** (1.40) | -0.13 (0.31) | -0.30** (0.14) |
| Θ | 0.15* (0.07) | 0.55** (0.23) | 0.24*** (0.07) | -0.12* (0.06) | -0.17* (0.08) | -0.27**** (0.03) |
| $ALTP_G$ | 6.68* (3.74) | -5.60 (6.39) | 7.41 (5.10) | 2.26 (2.87) | -5.13 (5.05) | -9.77**** (1.88) |
| Constant | 2.17*** (0.70) | 1.72 (1.00) | 2.75**** (0.67) | -0.17 (1.40) | 4.73** (2.03) | 6.25**** (0.82) |
| R-squared | 0.83 | 0.50 | 0.63 | 0.91 | 0.35 | 0.77 |
| Obs. (Clusters) | 45 (45) | 66 (10) | 111 (55) | 45 (45) | 66 (10) | 111 (55) |

Notes: Std. errors in parentheses, clustered at the issue level. *, **, ***, and **** denote significance at the 10% level, 5%, 1%, and 0.1% level, respectively. In pooled specifications, we include a dummy variable indicating whether the observation came from ANES data set.

Table 7: Prediction Errors of Representativeness-Based Model for GNH Data

| Representativeness Model: Truncation to Most Representative Types | | | | | | |
|---|---------|---------|---------|---------|---------|---------|
| | $d = 1$ | $d = 2$ | $d = 3$ | $d = 4$ | $d = 5$ | $d = 6$ |
| (a) Predicting Believed Typical Mean of Conservatives in GNH Data | | | | | | |
| Mean Squared Prediction Error | 3.02 | 1.05 | 0.54 | 0.30 | 0.27 | 0.48 |
| Mean Prediction Error | -1.36 | -0.62 | -0.28 | 0.040 | 0.33 | 0.56 |
| Rate of Underestimation | 3/45 | 7/45 | 13/45 | 25/45 | 35/45 | 41/45 |
| N | 45 | 45 | 45 | 45 | 45 | 45 |
| (b) Predicting Believed Typical Mean of Liberals in GNH Data | | | | | | |
| Mean Squared Prediction Error | 2.47 | 1.42 | 0.68 | 0.20 | 0.073 | 0.083 |
| Mean Prediction Error | 0.94 | 0.79 | 0.57 | 0.27 | 0.071 | -0.093 |
| Rate of Underestimation | 39/45 | 41/45 | 42/45 | 42/45 | 26/45 | 17/45 |
| N | 45 | 45 | 45 | 45 | 45 | 45 |

the true belief. We also present a simple counting metric: the fraction of observations for which the model underestimates the observed belief. We compute each of these three measures separately for each group, conservatives and liberals, for each value of d . The model with $d = 4$ or $d = 5$ produces smaller MSPE than the accurate beliefs benchmark ($d = 6$). Furthermore, in both cases, the errors are less systematic. While 41 of the 45 true means are less than the observed beliefs (indicating consistent underestimation), errors are more evenly distributed across over and underestimation for most of the truncation models. We do the same exercise for liberals in the GNH data in Panel (b).

In the main text, we compared this model to a model where beliefs are obtained by truncating to the d most likely types. Table 8 shows the MSPE, MPE, and rates of underestimation for the likelihood-based truncation model under different values of d . While the likelihood model produces smaller MSPE and MPE than the stereotype model for small values of d , for each group, the best representativeness-based model produces smaller MSPE errors than the best likelihood-based model. And, while the best representativeness-based model is a better predictor of observed beliefs than the accurate beliefs benchmark for both liberals and conservatives in terms of MSPE and MPE, the best likelihood-based model never beats the accurate beliefs benchmark in terms of MSPE. For conservatives, the best representativeness-based model outperforms the best likelihood-based model and the accurate beliefs benchmark for all metrics. For liberals, the best representativeness-based model

Table 8: Prediction Errors of Likelihood-Based Model for GNH Data

| Likelihood Model: Truncation to Most Likely Types | | | | | | |
|---|---------|---------|---------|---------|---------|---------|
| | $d = 1$ | $d = 2$ | $d = 3$ | $d = 4$ | $d = 5$ | $d = 6$ |
| (a) Predicting Believed Typical Mean of Conservatives in GNH Data | | | | | | |
| Mean Squared Prediction Error | 1.28 | 1.09 | 0.85 | 0.62 | 0.56 | 0.48 |
| Mean Prediction Error | 0.58 | 0.58 | 0.53 | 0.54 | 0.57 | 0.56 |
| Rate of Underestimation | 35/45 | 35/45 | 31/45 | 38/45 | 39/45 | 41/45 |
| N | 45 | 45 | 45 | 45 | 45 | 45 |
| (b) Predicting Believed Typical Mean of Liberals in GNH Data | | | | | | |
| Mean Squared Prediction Error | 1.17 | 0.61 | 0.31 | 0.18 | 0.12 | 0.083 |
| Mean Prediction Error | 0.52 | 0.37 | 0.21 | 0.077 | -0.018 | -0.093 |
| Rate of Underestimation | 32/45 | 32/45 | 32/45 | 29/45 | 24/45 | 17/45 |
| N | 45 | 45 | 45 | 45 | 45 | 45 |

Table 9: Prediction Errors of Representativeness-Based Model for ANES Data

| Representativeness Model: Truncation to Most Representative Types | | | | | | | |
|--|---------|---------|---------|---------|---------|---------|---------|
| | $d = 1$ | $d = 2$ | $d = 3$ | $d = 4$ | $d = 5$ | $d = 6$ | $d = 7$ |
| (a) Predicting Believed Typical Mean of Conservatives in ANES Data | | | | | | | |
| Mean Squared Prediction Error | 2.02 | 1.84 | 1.30 | 0.76 | 0.57 | 0.52 | 0.46 |
| Mean Prediction Error | -1.25 | -1.14 | -0.89 | -0.59 | -0.40 | -0.21 | 0.058 |
| Rate of Underestimation | 3/66 | 6/66 | 7/66 | 8/66 | 15/66 | 23/66 | 38/66 |
| N | 66 | 66 | 66 | 66 | 66 | 66 | 66 |
| (b) Predicting Believed Typical Mean of Liberals in ANES Data | | | | | | | |
| Mean Squared Prediction Error | 4.81 | 2.95 | 1.39 | 0.38 | 0.28 | 0.43 | 0.63 |
| Mean Prediction Error | 1.76 | 1.65 | 1.07 | 0.44 | 0.013 | -0.30 | -0.53 |
| Rate of Underestimation | 63/66 | 65/66 | 64/66 | 55/66 | 31/66 | 22/66 | 13/66 |
| N | 66 | 66 | 66 | 66 | 66 | 66 | 66 |

outperforms the best likelihood-based model and the accurate beliefs benchmark in terms of MSPE, with more mixed results for MPE.

Tables 9 and 10 repeat this exercise for the ANES data.

H.2 Beliefs of Conservatives and Liberals

In this section, we show that the predictions of our model hold both for beliefs held by Conservatives and beliefs held by Liberals. First, we document exaggeration in Table 11. In the GNH data, Liberals hold more exaggerated beliefs about both Conservatives and Liberals than Conservatives do. The pattern is different in the ANES data. Conservatives

Table 10: Prediction Errors of Likelihood-Based Model for ANES Data

| Likelihood Model: Truncation to Most Likely Types | | | | | | | |
|--|---------|---------|---------|---------|---------|---------|---------|
| | $d = 1$ | $d = 2$ | $d = 3$ | $d = 4$ | $d = 5$ | $d = 6$ | $d = 7$ |
| (a) Predicting Believed Typical Mean of Conservatives in ANES Data | | | | | | | |
| Mean Squared Prediction Error | 2.82 | 1.45 | 1.14 | 0.97 | 0.70 | 0.53 | 0.46 |
| Mean Prediction Error | -0.068 | -0.22 | -0.20 | -0.20 | -0.073 | 0.002 | 0.058 |
| Rate of Underestimation | 38/66 | 29/66 | 23/66 | 25/66 | 26/66 | 31/66 | 38/66 |
| N | 66 | 66 | 66 | 66 | 66 | 66 | 66 |
| (b) Predicting Believed Typical Mean of Liberals in ANES Data | | | | | | | |
| Mean Squared Prediction Error | 2.67 | 1.35 | 0.94 | 0.88 | 0.78 | 0.71 | 0.63 |
| Mean Prediction Error | -0.13 | -0.30 | -0.34 | -0.28 | -0.37 | -0.46 | -0.53 |
| Rate of Underestimation | 20/66 | 26/66 | 23/66 | 28/66 | 23/66 | 19/66 | 13/66 |
| N | 66 | 66 | 66 | 66 | 66 | 66 | 66 |

Table 11: Information about -G Predicts Beliefs about G, Conservatives versus Liberals

| Exaggeration of Beliefs about G | | | | |
|---------------------------------|-----------------------|------------------|-----------------------|------------------|
| | G = Conservatives | | G = Liberals | |
| | Held by Conservatives | Held by Liberals | Held by Conservatives | Held by Liberals |
| GNH | 0.35 | 0.71 | 0.03 | 0.21 |
| ANES | -0.11 | 0.18 | 0.78 | 0.36 |

in the ANES data have exaggerated beliefs about Liberals, but not about Conservatives. Liberals in the ANES data have exaggerated beliefs about both Liberals and Conservatives, with more exaggerated beliefs about Liberals than Conservatives. Given the differences across the two data sets, it is hard to draw general conclusions about whether beliefs are more exaggerated when predicting positions of the other group. In most cases, for both Liberals and Conservatives, reported beliefs are more extreme than the truth for both their own group and the other group.

Next, we test the prediction of Equation 5, asking whether we observe the same context-dependence for beliefs held by either group. In Table 12, we predict the believed mean of a group G from the true mean of the group G and the true mean of -G. Our model predicts that information about -G will be predictive of believed mean of G. The key here is whether this prediction holds independent of whether we are considering beliefs about G held by Conservatives or Liberals. Thus, we present two specifications side-by-side, one predicting beliefs held by Conservatives about a group G, and one predicting beliefs held by Liberals of

Table 12: Information about -G Predicts Beliefs about G, Conservatives versus Liberals

| OLS Predicting Believed Mean of G in Pooled Data | | | | |
|--|-----------------------|----------------------|-----------------------|----------------------|
| | G = Conservatives | | G = Liberals | |
| | Held by Conservatives | Held by Liberals | Held by Conservatives | Held by Liberals |
| True Mean Conservatives | 1.13**** (0.076) | 0.92**** (0.087) | -0.36**** (0.073) | -0.23**** (0.065) |
| True Mean Liberals | -0.55**** (0.118) | -0.65**** (0.149) | 0.68**** (0.140) | 0.81**** (0.147) |
| Constant | 1.46**** (0.279) | 2.87**** (0.305) | 2.14**** (0.260) | 1.16**** (0.282) |
| R-squared | 0.77 | 0.57 | 0.46 | 0.76 |
| Obs. (Clusters) | 111 (55) | 111 (55) | 111 (55) | 111 (55) |

Notes: Std. errors in parentheses, clustered at the issue level. *, **, ***, and **** denote significance at the 10% level, 5%, 1%, and 0.1% level, respectively. In pooled specifications, we include a dummy variable indicating whether the observation came from ANES data set.

that same group G. We see quite similar results when we explore beliefs held by Conservatives and beliefs held by Liberals. In particular, both sets of beliefs demonstrate the same strong evidence for context-dependence that we documented in the main text.

In Table 13, we test the predictions of Equations 6 and 7, asking whether Θ also has predictive power for beliefs held by Conservatives and Liberals. Again, we see that the results do not strongly depend on who holds the beliefs. In predicting the Conservatives' belief of the mean Conservative position or the Liberals' belief of the mean Conservative position, the average representativeness of tail positions has predictive power. Similarly, in predicting the Conservatives' belief of the mean Liberal position or the Liberals' belief of the mean Liberal position, the average representativeness of Liberal tail positions has a negative, but insignificant on beliefs.

Table 13: Average Representativeness of Tail Positions Predicts Beliefs, Conservatives versus Liberals

| OLS Predicting Believed Mean of G in Pooled Data | | | | |
|--|-----------------------|--------------------|-----------------------|--------------------|
| | G = Conservatives | | G = Liberals | |
| | Held by Conservatives | Held by Liberals | Held by Conservatives | Held by Liberals |
| True Mean of G | 0.71**** (0.09) | 0.42**** (0.10) | 0.36**** (0.10) | 0.58**** (0.10) |
| Θ_G | 0.20** (0.09) | 0.30**** (0.09) | -0.12 (0.10) | -0.07 (0.07) |
| Constant | 1.03**** (0.30) | 2.24**** (0.32) | 1.99**** (0.24) | 1.10**** (0.27) |
| R-squared | 0.72 | 0.50 | 0.32 | 0.75 |
| Obs. (Clusters) | 111 (55) | 111 (55) | 110 (54) | 110 (54) |

Notes: Std. errors in parentheses, clustered at the issue level. *, **, ***, and **** denote significance at the 10% level, 5%, 1%, and 0.1% level, respectively. In pooled specifications, we include a dummy variable indicating whether the observation came from ANES data set. One liberal observation is missing from the GNH data as there is no mass on stereotypical liberal positions for either group for one issue.

References for Online Appendix:

- Couch, James, and Jennifer Sigler. 2001. "Gender Perception of Professional Occupations." *Psychological Reports* 88 (3): 693 – 698.
- Decker, Wayne. 1986. "Occupation and Impressions: Stereotypes of Males and Females in Three Professions." *Social Behavior and Personality* 14 (1): 69–75.
- Hewstone, Miles, Manfred Hassebrauck, Andrea Wirth, and Michaela Waenke. 2000. "Pattern of Disconfirming Information and Processing Instructions as Determinants of Stereotype Change." *British Journal of Social Psychology* 39: 399 – 411.
- Kahneman, Daniel, and Shane Frederick. 2005. "A Model of Heuristic Judgment," in *The Cambridge Handbook of Thinking and Reasoning*, Keith Holyoak and Robert Morrison, eds. Cambridge, UK: Cambridge University Press.
- Lord, Charles, Lee Ross, and Mark Lepper. 1979. "Biased Assimilation and Attitude Polarization: The Effects of Prior Theories on Subsequently Considered Evidence." *Journal of*

Personality and Social Psychology 37 (11): 2098 – 2109.

Madon, Stephanie, Max Guyll, Kathy Aboufadel, Eulices Montiel, Alison Smith, Polly Palumbo, and Lee Jussim. 2001. “Ethnic and National Stereotypes: The Princeton Trilogy Revisited and Revised.” *Personality and Social Psychological Bulletin* 27 (8): 996 – 1010.

Nickerson, Raymond. 1998. “Confirmation Bias: a Ubiquitous Phenomenon in Many Guises.” *Review of General Psychology* 2 (2): 175 – 220.

Noori, Kamyar, and Allyson Weseley. 2011. “Beyond Credentials: The Effect of Physician Sex and Specialty on How Physicians Are Perceived.” *Current Psychology* 30: 275 – 283.

Rothbart, Myron. 1981. “Memory Processes and Social Beliefs.” In *Cognitive Processes in Stereotyping and Intergroup Behavior*, ed. DL Hamilton, Hillsdale, NJ: Erlbaum: 145 – 81.

Schneider, David. 2004. *The Psychology of Stereotyping*. New York, NY: The Guilford Press.

Schwartzstein, Joshua. 2014. “Selective Attention and Learning.” *Journal of the European Economic Association* 12 (6): 1423 – 1452.

Weber, Renée, and Jennifer Crocker. 1983. “Cognitive Processes in the Revision of Stereotypic Beliefs.” *Journal of Personality and Social Psychology* 45 (5): 961 – 977.