

Associative Judgment and Vector Space Semantics

Sudeep Bhatia
University of Pennsylvania

I study associative processing in high-level judgment using vector space semantic models. I find that semantic relatedness, as quantified by these models, is able to provide a good measure of the associations involved in judgment, and, in turn, predict responses in a large number of existing and novel judgment tasks. My results shed light on the representations underlying judgment, and highlight the close relationship between these representations and those at play in language and in the assessment of word meaning. In doing so, they show how one of the best-known and most studied theories in decision making research can be formalized to make quantitative a priori predictions, and how this theory can be rigorously tested on a wide range of natural language judgment problems.

Keywords: judgment and decision making, associative judgment, conjunction fallacy, distributional semantics, word vectors

Supplemental materials: <http://dx.doi.org/10.1037/rev0000047.supp>

Associations play a fundamental role in judgment. They process co-occurrence-based statistical regularities in a fast, automatic, and relatively effortless manner. Judgments relying on associations use these regularities to infer the suitability of a response, so that responses that are strongly associated with the content of the judgment question are the ones most likely to be selected by decision makers (Kahneman, 2003; Kahneman & Frederick, 2002; Morewedge & Kahneman, 2010; Sloman, 1996; also Evans, 2008 and Gilovich, Griffin, & Kahneman, 2002).

We currently do not have a computational or mathematical specification of associative judgment processes that is able to provide quantitative a priori predictions regarding the strength of association between a given question and response option, and thus provide quantitative a priori predictions regarding the judgments of decision makers. This is understandable. The types of problems used in judgment research can span a large domain of knowledge, and formalizing associative judgment processes so that they are able to encode this knowledge, specify associations, and make judgments almost universally, presents novel technical and theoretical hurdles for scholars of judgment and decision making.

An approach that is able to formalize these processes and quantitatively specify associations would, however, be of significant value. First, it would answer the call for more rigorous theories of associative judgment (Gigerenzer, 1996, 1998; Gigerenzer & Regier, 1996), and by doing so, make these theories easier

to test. Additionally, the predictions obtained from such an approach could be fit to both existing and novel judgment problems, and thus be used to determine not just whether associative processing describes the broad patterns underlying judgment, but also whether it can predict the specific response selections and probability assignments made by decision makers.

Such an approach would also be useful for studying responses to naturalistic judgment problems, that is, problems faced by decision makers in real-world settings. This would reveal both the external validity and the adaptive value of associative judgment. Additionally, such an approach would shed light on the mechanisms involved in representing information in high-level judgment tasks. Although there has been much work on the rules, strategies, and heuristics used by decision makers to weigh and aggregate information (Busemeyer, Pothos, Franco, & Trueblood, 2011; Dougherty, Gettys, & Ogden, 1999; Gigerenzer & Gaissmaier, 2011; Gilovich et al., 2002; Hammond & Stewart, 2001; Kahneman & Tversky, 1973; Reyna, Lloyd, & Brainerd, 2003; Shah & Oppenheimer, 2008; Tversky & Koehler, 1994; see also Weber & Johnson, 2009 and Hastie, 2001 for reviews), the question of what this information is has received little attention. It is clear that theories that do not address the issue of representation are unable to provide a complete account of the psychology of judgment.

In this article I examine an approach that makes explicit the associative relationships at play in judgment. As suggested in prior work (Kahneman, 2003; Kahneman & Frederick, 2002; Sloman, 1996), I assume that the questions in the judgment problems offered to decision makers activate the candidate response options based on the associations between the responses and the question. The relative activation of a response is used as a heuristic to judge the accuracy of the response, so that the response with the strongest association with the question in consideration is the one that is selected. Unlike prior work, however, I use vector space semantic models to specify the associations between a question and its candidate response options (Dhillon, Foster, & Ungar, 2011; Griffiths, Steyvers, & Tenenbaum, 2007; Jones & Mewhort, 2007;

Funding was received from the National Science Foundation Grant SES-1626825. This work was presented at meetings for the Cognitive Science Society (2016), the International Conference on Thinking (2016), the Psychonomic Society (2015), the Society for Judgment and Decision Making (2015), and the Basel workshop on memory and cognition (2015). Thanks to Nazli Bhatia for valuable feedback on earlier drafts of this article.

Correspondence concerning this article should be addressed to Sudeep Bhatia, Department of Psychology, University of Pennsylvania, 3720 Walnut Street, Philadelphia, PA 19104. E-mail: bhatiasu@sas.upenn.edu

Kwantes, 2005; Landauer & Dumais, 1997; Lund & Burgess, 1996; Mikolov et al., 2013; Pennington, Socher, & Manning, 2014). Vector space models are popular tools in computational linguistics and semantic memory research for studying word meaning, and have been shown to accurately predict phenomenon as diverse as semantic priming, synonym judgment, analogical judgment, and judgments of word association. The semantic relatedness between a pair of words corresponds to a form of association between the concepts corresponding to the words, and I suggest that semantic relatedness, as assessed by vector space models, can provide an accurate measure of the associations used in high-level judgment.

Associative Judgment

The representativeness heuristic is one of the best-known and most-studied theories in judgment and decision making research. In their classic 1974 article, Tversky and Kahneman described this heuristic as a way to answer questions of the following type: What is the probability that A belongs to/originates from/generates B? According to Tversky and Kahneman, decision makers do not consider probabilistic or logical relationships between A and B when answering these types of questions. Rather they make their judgments based on whether A is representative of, that is, similar to, B. Similarity is an important feature of cognition, and judgments using similarity can be made with relative ease. Indeed Tversky and Kahneman (1974) found that the representativeness heuristic could predict participant responses in a range of decision problems of the above type, including problems in which the heuristic generated an incorrect response (see also Kahneman & Tversky, 1973; Tversky & Kahneman, 1983).

In recent years the focus of decision making research has shifted from specifying individual heuristics to building more general frameworks within which heuristics and related psychological processes can be seen to operate. In this light, Kahneman and coauthors have suggested that judgments from representativeness are a product of an intuitive system that relies primarily on associations, and the response activations that they generate (Kahneman, 2003; Kahneman & Frederick, 2002; Morewedge & Kahneman, 2010). This system is relatively quick and effortless, but is unable to reason about abstract concepts and complex relationships, process logical structure, or use sophisticated decision rules. These properties, according to Kahneman et al., belong to a second, deliberative judgment system that operates in a controlled, but slow and effortful manner. Although decision makers can utilize both systems, the associative system is typically the first to engage and provide a response, and thus is more likely to influence judgment.

A dual-systems theory featuring both associative processes and rule-based processes has also been proposed by Sloman (1996; see also Barbey & Sloman, 2007). In addition to discussing the broad properties of these two systems and explaining how their interplay can organize key empirical findings, Sloman has suggested that these two systems utilize different forms of computation. The associative system can be seen as being instantiated in a connectionist network, which operates in parallel and draws inferences using co-occurrence based statistical regularities. In contrast, the deliberative system can be seen as being instantiated in a rule-

based architecture, which processes symbolic representations and abstracted logical relationships.

This article attempts to formally specify the associative judgment processes proposed in the work of Kahneman, Sloman, and others (see Evans, 2008 for a review). The approach used in this article is motivated by Sloman's claim that associative judgment processes are best understood as utilizing co-occurrence-based statistical regularities, and by Kahneman et al.'s claim that the response options activated and made accessible by these associative processes are the ones that are most likely to be selected by decision makers.

To see how associative processes of this type make judgments, let us consider the Linda problem (Tversky & Kahneman, 1983). In this problem decision makers are given the following description: *Linda is 31 years old, single, outspoken and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in antinuclear demonstrations.* They are then asked whether she is more likely to be a bank teller or a feminist bank teller. Decision makers typically believe that Linda is more likely to be a feminist bank teller than a bank teller, despite the fact that the set of feminist bank tellers is a subset of the set of bank tellers, making it impossible that Linda is a feminist bank teller but not a bank teller. This is a conjunction fallacy.

The biases generated by associative judgment processes provide a compelling explanation for this important finding. Feminists are often intellectually inclined, and concerned with discrimination and social justice, and thus instances of feminism frequently occur with an interest in philosophy and concern for social justice. Because of this co-occurrence, and the tendency of mind to learn co-occurrence-based statistical relationships, the description of Linda is associated with that of a feminist. As a result, activating mental representations corresponding to Linda's description activates mental representations of feminism and feminists. Finally, as strength of activation determines responses, decision makers are more likely to think that Linda is a feminist bank teller compared to just a bank teller, generating the conjunction fallacy.

Vector Space Semantic Models

A feminist bank teller is considered to be associated with Linda, because people with Linda's description are often feminists. However, how can we formally specify the associative connection between a feminist bank teller and Linda? Or, more generally, how can we specify the strength of associations between a given question and its candidate response options, so as to quantitatively predict, a priori, the probability assignments and response selections of decision makers to that question? Doing so is necessary to fully understand the predictions of theories of associative judgment, and in turn to rigorously test these theories on experimental and real-world judgment problems.

Although research on associative judgment is silent on this issue, the problem of understanding associations and, more generally, co-occurrence-based statistical regularities, has received considerable attention in semantic memory research and computational linguistics. Here scholars have used these types of regularities not to uncover the associations at play in judgment, but rather to uncover the representations involved in processing word meaning and the semantic relatedness between sets of words. The

key insight underlying this research is that words that have similar distributions in language have similar meanings, so that the semantic relationships between words can be determined by studying which words co-occur with each other. This approach to understanding word meaning has a long history in psychology (Firth, 1957; Harris, 1954), but has recently received increased attention with the successes of vector space semantic models (Dhillon et al., 2011; Griffiths et al., 2007; Jones & Mewhort, 2007; Kwantes, 2005; Landauer & Dumais, 1997; Lund & Burgess, 1996; Mikolov et al., 2013; Pennington et al., 2014). These models characterize each word in their vocabulary as a vector in a multidimensional space. The proximity between the vectors of two words corresponds to the semantic relatedness of the words, so that words that are distributionally similar to each other are located near each other in the multidimensional space, and are considered to be related to each other.

Vector space models are trained on very large natural language text corpora, and have large vocabularies, which can be used to make predictions regarding judgments of word similarity and synonymy, word analogy, the strength of word priming, word categorization, word association, reading times, recall times, and related psycholinguistic phenomena, for nearly all the words commonly used in a given language. The predictions of these models have been shown to be highly accurate, suggesting that the representations recovered by these models provide a good characterization of the representations underlying semantic processing in language (see Bullinaria & Levy, 2007 or Jones, Willits, & Dennis, 2015 for a review). For this reason, these models are also popular in machine learning and artificial intelligence, particularly in applications related to computer processing of natural language (Turney & Pantel, 2010).

Vector space models are typically only applied to language-based tasks. It is possible, however, that these models can also predict responses in more complex domains, such as high-level judgment. The fact that these models are able to account for phenomena such as priming, implies that the co-occurrence-based statistical regularities captured by these models influence automatic processing, and are able to generate the types of activation-based associative response biases that are commonly assumed to underlie judgment (Kahneman, 2003; Kahneman & Frederick, 2002; Sloman, 1996). This is also suggested by the work of Paperno, Marelli, Tentori, and Baroni (2014), who find that distributional measures of word association correlate very strongly with explicit probability judgments of word co-occurrence. Indeed, most recently, Kintsch (2014) has suggested that latent semantic analysis (Landauer & Dumais, 1997), a very prominent vector space model, could be applied to certain types of judgment by combining it with Busemeyer et al.'s (2011) quantum judgment theory.

To examine how vector space models could be applied to associative judgment, realize that many of the judgment problems in decision making research, such as the Linda problem described above, involve a question followed by a set of feasible response options. Here both the question and the various responses are pieces of text composed of words, and vector representations of the individual words in these texts can be used to generate aggregate vector representations for the question and responses. The proximity between these aggregate vector representations can, in turn, be used to determine the strength of association between the

question and the various responses, and subsequently predict the probabilities assigned to the response options by decision makers and the response options chosen by decision makers.

In this article I consider three state-of-the-art techniques for building vector representations, and examine whether pretrained vector representations based on these techniques are able to predict the associations at play in high-level judgment. Note that I am not attempting a comparison between these representations or techniques. Rather my goal is to use these different representations to establish the robustness of the proposed approach.

The first technique for building vector representations that I consider involves vector representations obtained using the continuous bag-of-words (CBOW) and skip-gram methods of Mikolov et al. (2013a, 2013b). This approach relies on a recurrent neural network that, for the CBOW method, predicts words using other words in their immediate context, and for the skip-gram method, attempts to do the inverse of this. The word vectors I use for this method are vectors released by Google Research through the Word2Vec tool. These vectors were trained on a corpus of Google News articles with over 100 billion words tokens, and have a vocabulary of 300 million words and phrases, with each word or phrase being defined on 300 dimensions.

I also use vector representations trained with the Eigenwords technique, which relies on canonical correlation analysis (CCA; Dhillon et al., 2011; Dhillon, Foster, & Ungar, 2015; Dhillon, Rodu, Lu, Foster, & Ungar, 2013). CCA is a dimensionality reduction method similar to the singular value decomposition that is used in latent semantic analysis (Landauer & Dumais, 1997). In this article I use a set of pretrained Eigenwords vectors released publicly by Dhillon et al. These vectors have been trained on the English Gigaword corpus, which is a comprehensive archive of newswire text data. The Eigenwords vectors have a vocabulary of 300,000 words with each word defined on 200 dimensions.

The final technique involves the Global Vectors for Word Representation (GloVe) model (Pennington et al., 2014), which performs a dimensionality reduction on word co-occurrence matrices, emphasizing the use of the ratios of word-word co-occurrence probabilities. In this article I use publicly available GloVe vectors obtained from Pennington et al. These vectors have been trained on a 6 billion word corpus combining English language Wikipedia with the English Gigaword corpus. They have a vocabulary of 400,000 words and 300 dimensions.

Now, the question in a given judgment problem can be represented in a bag-of-words format, as a collection of its component words. Vector space models specify each of these words as a vector w_i . It is possible to generate an aggregate representation of the question by taking the average of these vectors, weighted by the frequency of their corresponding words in the question. Thus the vector corresponding to the question, q , can be written as

$$q = \frac{\sum_i n_i w_i}{\sum_i n_i},$$

where n_i is the number of times word i occurs in the question. We can use the same method to build a vector representation of the response r , and in turn specify the association between the question and the response based on the distance between q and r . I use cosine similarity to specify distance, so that the association between q and r is $A(q,r) = q \cdot r / (\|q\| \cdot \|r\|)$. In the subsequent sections I will be considering problems in which a set of feasible responses are given to the decision maker. I will normalize the

associations for the responses in a given question using a softmax rule, so that for a question with feasible responses $r_1, r_2 \dots r_N$, the normalized association of response i with the question is written as $\tilde{A}(q, r_i | r_1, r_2 \dots r_N) = \frac{e^{A(q, r_i)}}{\sum_k e^{A(q, r_k)}}$. This parameter free method ensures that the probabilities predicted by this approach lie between 0 and 1, and sum to 1. Additional details regarding this approach are provided in the online supplemental materials.

Before continuing let us summarize why it is reasonable to assume that the vector space representations generated by Word2Vec, Eigenwords and GloVe, and representations obtained through related techniques, capture the associations at play in judgment. First, as discussed above, both theories of associative judgment and distributional models of semantics (such as the vector space models considered here) rely on co-occurrence-based statistical regularities to specify the various relationships between judgment objects or words, and in turn the influence of these relationships on automatic, activation-based response biases. This suggests that the associations uncovered by one approach can be used to understand the associations used by the other. Second, the processes involved in judgment apply primarily to information stored in the minds of decision makers, and vector space models provide one of the most powerful tools for understanding how this information is learnt and how it is represented. As such, knowledge about the objects at play in judgment can be reasonably assumed to stem from knowledge about the meanings of the words that are typically used to describe these objects. Finally, because of recent advances in machine learning, these models can be trained on very large natural language data sets, and as a result, possess expansive vocabularies. Models of associative judgment, once equipped with vector representations for these vocabularies, can subsequently be used to predict judgments for a very large number of commonly used natural language problems in judgment and decision making research. This would not be possible with other less computationally tractable approaches to specifying knowledge representations, which would have limited vocabularies and limited applicability. Thus, the use of vector space models in judgment is not only desirable from a theoretical perspective, but from a practical perspective as well.

Predicting Fallacies

In the previous section, I have outlined a method for formalizing the theory of associative judgment and for testing it rigorously on experimental data. In this section I perform such tests on the judgment fallacies commonly attributed to associative processing.

Linda Problem

The best known example of fallacious judgment generated by the use of associations involves the Linda problem (Tversky & Kahneman, 1983). Using the vector space representations described above, we can specify each of the words in the question in the Linda problem as a multidimensional vector, and subsequently average the set of words in the question to obtain a single vector describing Linda. I will refer to this vector as q_L . We can do the same with the two response options, *bank teller* and *feminist bank teller*, to obtain response vectors r_{BT} and r_{FBT} . We can subsequently write the associations between the two response options

and the question as $A(q_L, r_{BT})$ and $A(q_L, r_{FBT})$, and predict the options most likely to be chosen by decision makers based on whether $A(q_L, r_{BT}) > A(q_L, r_{FBT})$ or $A(q_L, r_{BT}) < A(q_L, r_{FBT})$.

I find that all three of the representations specify a stronger association between a feminist bank teller and Linda, than that between a bank teller and Linda, that is, represent q_L, r_{BT} , and r_{FBT} in a manner such that $A(q_L, r_{FBT}) > A(q_L, r_{BT})$. This happens because all of these representations encode strong relationships between words like *outspoken*, *philosophy*, and *justice*, which make up the description of Linda, and the word *feminist*, which is part of the response *feminist bank teller*. In contrast, words such as *bank* and *teller* do not have strong relationships with the words in Linda's description. Subsequently, the vectors for the words in Linda's description are closer to the word *feminist* than they are to the words *bank* and *teller*, and the aggregate vectors q_L and r_{FBT} are closer together compared with q_L and r_{BT} .

Although decision makers display a conjunction fallacy between the responses *bank teller* and *feminist bank teller*, they do not do so for the responses *feminist* and *feminist bank teller* (Tversky & Kahneman, 1983). The proposed approach is able to predict this because of the fact that the vector corresponding to the response *feminist*, r_F , is relatively close to the words in the description of Linda. In contrast, the words *bank* and *teller*, and subsequently the aggregate vector for *feminist bank teller*, r_{FBT} , is more distant from the words in Linda's description. Ultimately, responses that are composed primarily of words that are semantically related to the words in Linda's description will be more strongly associated with Linda, than other responses whose component words also include words that are not semantically related to Linda.

Figure 1 (left panel) illustrates this insight in a simple two dimensional space. It shows hypothesized vectors corresponding to the aggregate representation of Linda, as well as the aggregate response vectors corresponding to *feminist*, *feminist bank teller*, and *bank teller*. As outlined above, these aggregate vectors are obtained by averaging the vectors for their component words. Because r_F is closest to q_L , and r_{BT} is furthest from q_L , r_{FBT} lies between r_F and r_{BT} , and subsequently ends up being the second closest response to q_L . This generates $A(q_L, r_F) > A(q_L, r_{FBT}) > A(q_L, r_{BT})$, which corresponds to a conjunction fallacy between bank teller and feminist bank teller, but not one between feminist and feminist bank teller.

The associations generated by each of the vector representations for the responses in the Linda problem are shown in Figure 2. To facilitate comparisons between questions, Figure 2 displays nor-

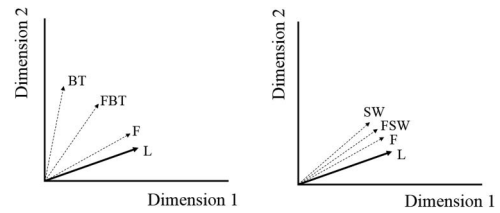


Figure 1. Example of vector representations for the responses (bank teller, BT; feminist bank teller, FBT; social worker, SW; feminist social worker, FSW; and feminist, F) to the Linda problem, as well as the vector corresponding to Linda (L). The association of a response with Linda is judged by its cosine similarity with the Linda vector.

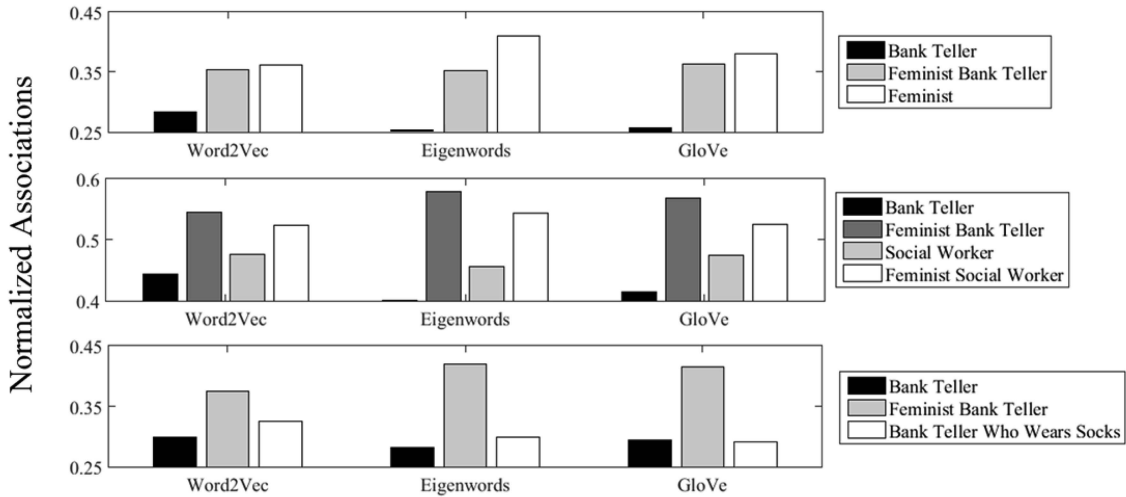


Figure 2. Normalized associations between Linda's description and various response options, generated by the three sets of vector representations.

malized associations. As shown in this figure, all three of these representations are able to generate the pattern of responses observed in experiments, with $\tilde{A}(q_L, r_F) > \tilde{A}(q_L, r_{FBT}) > \tilde{A}(q_L, r_{BT})$. Note that these associations do not vary significantly if the wording in the response *feminist bank teller* is replaced with *bank teller who is active in the feminist movement*.

In addition to the Linda problem, Tversky and Kahneman (1983) also discussed the problem of Bill, who was described in the following manner: *Bill is 34 years old. He is intelligent, but unimaginative, compulsive, and generally lifeless. In school, he was strong in mathematics but weak in social studies and humanities.* Participants were asked to judge the probabilities of the following statements: *Bill plays jazz for a hobby* and *Bill is an accountant who plays jazz for a hobby*, and, as with the Linda problem, these participants generated the conjunction fallacy by attaching higher probabilities to the second response. I also find that all three of the vector representations attach a stronger association to the response involving Bill being an accountant who plays jazz as a hobby, than they do to Bill playing jazz as a hobby. These results are summarized in Figure 3.

A related problem involves Danielle, described in the following manner: *Sensitive and introspective. In high school she wrote poetry secretly. Did her military service as a teacher. Though beautiful, she has little social life, since she prefers to spend her time reading quietly at home rather than partying.* Participants were asked whether Danielle is more likely to study literature or study the humanities (Bar-Hillel & Neter, 1993). Here participants typically believed that Danielle is more likely to study literature, despite the fact that the set of literature students is contained in the set of humanities students. As with the conjunction fallacy, this problem involves a setting in which the use of associations leads to participants ignoring category membership relationships. However, it is not a conjunction that is believed to be more probable than its constituents, but rather an implicit disjunction that is believed to be less probable than its constituents. Additionally, this disjunction fallacy does not emerge for all categories: decision makers are typically able to avoid this fallacy when asked whether Danielle is more likely to study physics or study the natural sciences. Again, the three sets of vector representations are able to predict these relationships, indicating that they provide a good

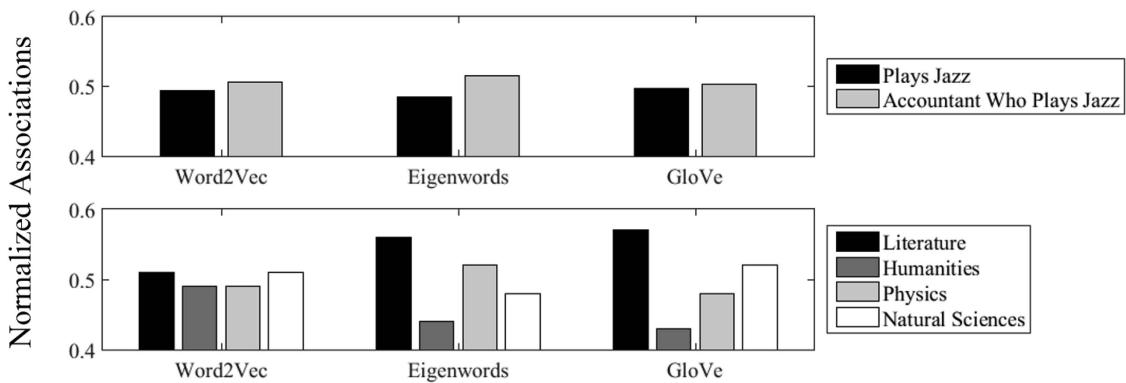


Figure 3. Normalized associations of various responses with the descriptions of Bill (top panel) and Danielle (bottom panel) generated by the three sets of vector representations.

account of different variants of the conjunction fallacy. These results are also summarized in Figure 3.

Typicality

After Tversky and Kahneman's seminal work there have been a number of researchers who have attempted to outline moderators of the conjunction fallacy. One early attempt at this involves variants of the Linda problem suggested by Shafir, Smith, and Osherson (1990). Shafir et al. propose that decision makers use the degree to which responses are typical of the descriptions in the question, to judge response probabilities. This leads to the prediction that the conjunction fallacy should be weaker when a conjunction like *feminist bank teller*, whose constituents are incompatible, is replaced with a compatible conjunction such as *feminist social worker*.

Shafir et al. (1990) verified this prediction experimentally, and I find that these results are also obtained using the proposed approach. Particularly, all three of the vector representations predict that the relative association between Linda's description and feminist social worker compared to Linda's description and social worker is smaller than that between Linda's description and feminist bank teller compared with Linda's description and bank teller. The reason for this can again be understood in terms of vector arithmetic. The vector for *feminist*, r_F , and the vector for *social worker*, r_{SW} , are both close to each another and close to the vector for Linda's description, q_L . Thus, averaging the vector for *social worker* with that for *feminist* to generate a vector for *feminist social worker*, r_{FSW} , does not lead to a large change in proximity to q_L compared to using just r_{SW} alone. Subsequently $A(q_L, r_{FSW}) - A(q_L, r_{SW})$ (and equivalently $\tilde{A}(q_L, r_{FSW}) - \tilde{A}(q_L, r_{SW})$) is relatively small. This is in contrast to the averaging that is done for *feminist bank teller*, for which r_F and r_{BT} are far apart, causing $A(q_L, r_{FBT}) - A(q_L, r_{BT})$ to be fairly large. This insight is again depicted in Figure 1 (right panel). Additionally, the associations for this example are summarized in the second panel of Figure 2.

Overall, Shafir et al. (1990) presented participants with multiple descriptions of individuals with both compatible and incompatible conjunctions. In each of these problems participants were asked to attach a probability to the described individual being in the two categories corresponding to the responses. The proposed approach is able to make precise quantitative predictions regarding the associations between any given natural language question and a

response, and thus can be used to predict the observed responses in Shafir et al.'s data.

There are a total of 56 probability assignments (obtained by averaging participant data) over the 14 different types of problems in Shafir et al.'s data, and I predict these average probability assignments using normalized associations for the questions and their corresponding responses obtained from the three vector space representations. The predictions are fit using a linear regression, permitting random intercepts for different questions. I find strong positive relationships between the normalized associations generated by all three representations and the responses of participants ($p < .05$ for Word2Vec and Eigenwords, and $p < .01$ for GloVe), indicating that the vector representations do not only describe the qualitative patterns observed in Shafir et al.'s data but are also able to match the specific probability assignments observed in this data. Additional details regarding the data, analysis, and results are provided in Table 1 and in the online supplemental materials.

Averaging

The changes in relative association outlined above are a product of the way in which the vector space approach averages the vectors of the words in the responses. Essentially, averages of two extremely proximate words or phrases (such as *feminist* and *social worker*) will be closer to the words themselves compared to averages of two distant words (such as *feminist* and *bank teller*), causing to responses composed of proximate words to have similar distances to the question. The insight that averaging plays a critical role in explaining the strength of the conjunction fallacy is not new to this article. A number of important explanations for this fallacy have involved averaging the probabilities for the individual components of the conjunction to generate an overall probability for the conjunction (Fantino, Kulik, Stolarz-Fantino, & Wright, 1997; Nilsson, Winman, Juslin, & Hansson, 2009; Yates & Carlson, 1986). The use of this averaging is incompatible with the type of probability aggregation permitted by probability theory: The probability of a conjunction can never be greater than the probability of its least probable component, and thus can never actually be the average of its components.

Gavanski and Roskos-Ewoldsen (1991) performed an early experiment examining this averaging hypothesis in detail. They asked participants to evaluate responses with high and low probabilities as well as conjunctions of these responses. If averaging

Table 1

Summary of Model Fits of the Three Vector Representations to the Four Existing Datasets and One Novel Dataset, Regarding Judgment Fallacies

Data	Word2Vec			Eigenwords			GloVe		
	β	z	R^2	β	z	R^2	β	z	R^2
Shafir et al. (1990)	2.21	2.24**	.08	1.87	2.28**	.09	2.01	2.94***	.14
Gavanski and Roskos-Ewoldsen (1991)	9.04	2.56**	.13	3.64	1.74*	.08	6.32	2.81***	.16
Tentori et al. (2013)	4.31	2.23**	.33	3.19	1.92*	.27	3.13	2.47**	.38
Study 1	1.05	20.03***	.55	1.04	19.41***	.54	1.06	20.20***	.56
Fischhoff and Bar-Hillel (1984)	7.55	4.56***	.27	2.58	1.35	.04	4.42	3.86***	.21
Average R^2			.27			.20			.29

Note. The fits involve random-effects linear regressions to predict the responses of participants using the normalized associations generated by the vector representations.

* $p < .1$. ** $p < .05$. *** $p < .01$.

plays a role in the conjunction fallacy, then we should expect probabilities assigned to conjunctions to lie between the probabilities assigned to their various components, so that if the components are both individually assigned high probabilities (as would be the case for the responses *feminist* and *avid reader* to the Linda problem) then the conjunct (in this case, *an avid reader who is a feminist*) should too. Likewise, if the components both have low probabilities (e.g., *bank teller* and *fashion conscious*), the conjunct (*bank teller who is fashion conscious*) should too, and if one of the components has a high probability and the other has a low probability (e.g., *feminist* and *bank teller*), the conjunct (*feminist bank teller*) should have a moderate probability. Gavanski and Rosko-Ewoldsen found that these predictions held in their dataset, indicating that averaging plays a role in ascribing probabilities to conjunctions.

As discussed above, the proposed approach is able to predict this qualitative pattern of behavior. This approach essentially averages the vectors of the components of a conjunction to generate a vector for the conjunction, and thus the response probabilities for the conjunction typically lie between the response probabilities of the components. However, again we can do more than just mimic qualitative patterns; we can also try to quantitatively predict the responses in Gavanski and Rosko-Ewoldsen's (1991) dataset. I did this for the Linda-type problems in this dataset, which involve four different target descriptions, with each description having seven different unique responses (this data also had what Gavanski and Rosko-Ewoldsen refer to as mixed problems and probability combination problems, which I did not attempt to fit). For each of the descriptions and responses I used the three vector representations to generate normalized associations, and fit these normalized associations to the probability assignments of participants, with a linear regression, with random intercepts on the question level. I found that the associations generated the three vector representations all had positive relationships with the responses of participants ($p < .05$ for Word2Vec, $p < .10$ for Eigenwords, and $p < .01$ for GloVe). Additional details regarding the data, analysis, and results are provided in Table 1 and in the online supplemental materials.

Confirmation

Gavanski and Rosko-Ewoldsen (1991) provide evidence that people average the probabilities of the components of conjunctions to judge the overall probabilities of the conjunctions. This mechanism, which the proposed approach is often able to mimic because of its assumptions regarding vector averaging, is a valuable explanation for why the probability of Linda being a feminist bank teller lies between the probability of Linda being a bank teller and Linda being a feminist. However, simple averaging of objective probability by itself cannot provide a full account of the conjunction fallacy. Consider for example, the following response to the Linda problem, proposed by Tentori, Crupi, and Russo (2013): *Linda is a bank teller who wears socks*. If decision makers judge the probability of this conjunction by averaging the objective probabilities of the conjuncts then it must be the case that they assign a higher probability to the above response than they do to Linda being a feminist bank teller. This is because nearly everyone wears socks, and Linda, in turn, is more likely to wear socks than to be a feminist. Subsequently the average of the probabilities

assigned to Linda being a bank teller and Linda wearing socks must be higher than the average of the probabilities assigned to Linda being a feminist bank teller. Tentori et al. (2013), however, find that decision maker are more likely to select responses that appear to be confirmed or supported by the description in the judgment problem (responses such as *feminist bank teller*), rather than responses whose components have higher objective probabilities (responses such as *bank teller who wears socks*). As with the variants of the Linda problem discussed above, the proposed approach is also able to generate this pattern, giving stronger associations between Linda's description and *feminist bank teller*, compared with *bank teller who wears socks*. These associations are summarized in the third panel of Figure 2.

Tentori et al. (2013) examined the predictions of their confirmation-theoretic account across a number of experiments (also see Crupi, Fitelson, & Tentori, 2008). In their second experiment they considered four variants of the Linda problem, with each variant involving three possible response options: one isolated response (e.g., *bank teller*), one conjunction of this response with a probable response (e.g., *bank teller who wears socks*), and one conjunct of this response with a confirmed response (e.g., *feminist bank teller*). In all four of these problems they found that the conjunct involving the confirmed response is more likely to be chosen by participants than the conjunct involving the probable response.

I attempted to fit this data quantitatively by using the vector representations to generate normalized associations for the three responses for each of the four descriptions (other experiments in Tentori et al. involved standalone events and negated events, which I did not attempt to fit). These normalized associations were then fit to the response probabilities of participants, with a linear regression, with random intercepts on the question level. Although the dataset had only 12 responses, the associations generated by the vector space representations all had positive relationships with these responses, as is shown in Table 1 ($p < .05$ for Word2Vec and GloVe, and $p < .10$ for Eigenwords). Additional details regarding the data, analysis, and results are provided in Table 1 and in the online supplemental materials.

Participant-Generated Problems

The above sections involve a small number of problems that have been generated by experimenters. A rigorous quantitative test of our ability to predict behavior in conjunction fallacy problems would also benefit from a larger set of problems generated by participants themselves. For this purpose I conducted a novel experiment (Study 1) on the conjunction fallacy using participant-generated judgment problems. Particularly, in this experiment, I asked 50 Amazon Mechanical Turk (MTurk) participants to list a hobby or an interest, and in turn specify three adjectives that describe someone with that hobby, one occupation that someone with that hobby would be highly likely to have, and one occupation that someone with that hobby would be highly unlikely to have. Using these adjectives as person-level descriptions in the judgment question, and the hobbies, job occupations, and their conjunctions, as the possible response options, I was able to generate a large number of Linda-type conjunction fallacy problems. These problems were then administered to a second pool of 300 MTurk

participants. There were a total of 82 different problems and participants in the second pool were given only one problem each.

The responses of these participants were fit using the methods described above. Particularly, I used the three vector space representations to obtain normalized associations for each of the responses to the 82 judgment problems. These were fit to the average response proportions of participants, using a linear regression, with random intercepts on the problem level. I found that all vector representations had a positive relationship with the responses of participants ($p < .01$ for all three representations), indicating that the proposed approach is not only able to quantitatively describe participant behavior in the conjunction fallacy for experimenter generated problems, but is also able to do so when the problems themselves are generated by participants. Additional details regarding the experimental methods, analysis, and results are provided in [Table 1](#) and in the online supplemental materials.

Base Rates

Associative judgment processes are not only responsible for the conjunction fallacy. They have also been seen as a cause of the neglect of base rates, which is the finding that decision makers often place too little weight on the prior probabilities of the various response options. [Kahneman and Tversky \(1973\)](#) illustrated this effect by asking participants to judge the likely job occupations of hypothetical people with different descriptions, while also telling participants about the distribution of the job occupations in the population from which the person was drawn. They found that people typically ignore these distributions, so that changing base rates does not alter probability judgments. For example, an individual, Jack, was described as *a conservative and careful father who likes carpentry and mathematics puzzles and does not have any interest in political or social issues*. Participants were also either told that Jack was one of 30 engineers in a population of 100 engineers and lawyers or one of 70 engineers in a population of 100 engineers and lawyers. They were then asked to assign a probability to Jack being an engineer or a lawyer. [Kahneman and Tversky](#) found virtually no difference in these probability assignments across the base rate conditions.

These effects were further tested in [Fischhoff and Bar-Hillel \(1984\)](#), who gave participants multiple judgment problems of this type. If vector space models can adequately specify the associations at play in judgment we should expect these models to generate good fits to not only the conjunction fallacy data sets considered thus far, but also the dataset in [Fischhoff and Bar-Hillel \(1984\)](#). I tested this using the first experiment presented by [Fischhoff and Bar-Hillel](#). This experiment involved two sets of descriptions under two conditions which varied whether the base rates were 70 or 30% for the given response category. In total, this led to 60 different problems of the type initially used by [Kahneman and Tversky \(1973\)](#). I obtained normalized associations for each of the responses to these problems using the three vector representations, and fit these normalized associations to the response proportions of participants with a linear regression, controlling for the base rates of the response category in consideration. I found that, once again, all three of the vector representations generated positive relationships between normalized associations and the responses of [Fischhoff and Bar-Hillel's](#) participants. Although this time only the relationships generated by Word2Vec and GloVe

reached statistical significance ($p < .01$ for these two representations, $p > .20$ for Eigenwords) this nonetheless illustrates the unique power of the proposed approach. It is able to provide quantitative predictions regarding existing experimental data not only for the conjunction fallacy, but also for related fallacies such as base rate neglect. Additional details regarding the data, analysis, and results are provided in [Table 1](#) and in the online supplemental materials.

Alternate Models of Judgment Bias

The above results show that the proposed approach presents a powerful technique for formalizing associative judgment and predicting judgment errors such as the conjunction fallacy. Of course associative judgment is not the only theory predicting these behaviors, and so it is useful to consider the relationship between the proposed approach and various models proposed in prior work. As discussed above, one important set of prior models involves averaging rules ([Fantino et al., 1997](#); [Nilsson et al., 2009](#); [Yates & Carlson, 1986](#)). These rules predict that probabilities attached to conjunctions and disjunctions are a weighted aggregate of the probabilities attached to their individual components, though the specific weights can vary based on the logical connective used (for this reason, these models are also sometimes referred to as configural weighting rules). Many of the predictions of these rules are shared by the proposed approach, which, through vector averaging, often leads to probability judgments resembling weighted averaging. Of course, there are some key differences as well: Averaging models assume that individuals have some probabilistic beliefs regarding individual events, and attempt to model how these probabilities are aggregated in to conjunctions. In contrast, the proposed approach attempts primarily to address how the underlying probability estimates are formed. The averaging-type behavior of our model is a merely a useful byproduct of the vector aggregation rules that are assumed to be at play in generating representations for complex propositions.

A second prominent theory of probability judgment biases involves quantum probability ([Busemeyer et al., 2011](#); [Franco, 2009](#); [Trueblood & Busemeyer, 2011](#)). According to this approach, beliefs are vectors and events are subspaces in a multidimensional space, and probability estimates are formed by projecting belief vectors onto the event subspace. Conjunction errors may arise when events are incompatible, and the strength of these errors depends on the order in which the different events are evaluated. This is quite a powerful theory, and is able to describe a large number of departures from probability theory observed in judgment. However, as with averaging models, the focus of quantum judgment theory is on proposing rules for combining conjunctions, disjunctions, and other logical connectives. Unlike the proposed approach, the underlying representations to which quantum judgment rules are applied are typically not specified in an a priori manner. Despite these differences, the use of vector calculations in quantum judgment does resemble the vector-based analysis outlined above, suggesting that the proposed approach could be extended using insights from quantum theory. I consider this possibility in the discussion section of this article.

Another important theory closely related to the proposed approach involves confirmation, as in the work of [Tentori et al. \(2013\)](#); see also [Crupi et al., 2008](#)). As discussed above, this work

suggests that decision makers judge the probability of an event or outcome based on the degree to which evidence confirms (i.e., increases the posterior relative to prior probability of) the event. Although the notion of confirmation studied by Tentori et al. (2013) comes from a very different research tradition compared to the word co-occurrence-based associative relationships I examine in this article, the two constructions are closely related. Particularly, if a hypothesis is independent of a piece of evidence, then the evidence does not confirm or disconfirm the hypothesis. Independence also implies that the variables that correspond to the evidence and the hypothesis are uncorrelated, which means that these variables do not co-occur systematically. For this reason we would expect a response that is confirmed by the description in the question to also be strongly associated with this description. Indeed, the relationship between association (as modeled in this article) and confirmation (as modeled in the work of Tentori et al.) may be more intimate than this: Association can be seen as representing the way in which confirmation is assessed in natural language judgment problems (see Paperno et al., 2014 for a related point), and, given the primacy of associative processing, can explain why people use confirmation (rather than probability) in making Linda-type judgments. In turn, confirmation can specify the epistemic and logical properties of associative processing with rigor, and can more generally be used to study judgments in settings in which linguistic associations are not directly applicable.

Yet another explanation for the conjunction fallacy relies on noisy recall. This account, formalized by Costello and Watts (2014), suggests that decision makers have representations of events in their memory. Probability judgments for these events are formed by recalling these events with some error, which can generate conjunction fallacies (Costello, 2009; Erev, Wallsten, & Budescu, 1994; Hilbert, 2012 provide related noise-based accounts of judgment bias). Although the noisy recall model approach is able to explain certain patterns observed with regards to the conjunction fallacy in human data, observed rates of conjunction fallacies are often higher than rates that can be predicted by this approach (see Crupi & Tentori, 2016 and Nilsson, Juslin, & Winman, 2016 for a critique). Additionally, unlike the proposed approach, the noisy recall model is unable to specify just what the event representations are, and subsequently unable to predict, a priori, the probability assignments of individuals. However, note that the vector space approach proposed in this article does not currently model noise in the associative judgment process, suggesting that some of the assumptions of the noisy recall model could be extended to the proposed approach to improve its predictions.

The noisy recall model is also closely related to Minerva-DM (Dougherty et al., 1999), which has a key memory component. Unlike the noisy recall model, which counts up events in memory, Minerva-DM uses similarity with a probe to make judgments. Noise is most likely to bias these judgments when one of the conjuncts is very similar to the probe. Minerva-DM is a very powerful theory of judgment, and is able to explain a number of key findings in the literature. More important, its use of similarity-based memory processes suggests a close relationship with the proposed account of the conjunction fallacy: Both theories rely critically on the strength of retrieved representations. However, again, as with many of the models outlined above, the focus of Minerva-DM is not on specifying what these memories and rep-

resentations are. Rather this work attempts to formalize the processes involved in memory-based judgment. This is contrast to the proposed approach, which comes equipped with the actual associations necessary to make predictions in a given (natural language) problem.

Ultimately, the study of the conjunction fallacy and related judgment biases has a rich theoretical and experimental history, with many documented findings, divergent task and problem representations, and competing explanations. The vast scope of this research suggests that the proposed approach (and associative processing more generally) may not be able to provide a conclusive account of all of the subtleties of human judgment. For example, the use of word vector averaging is a common way to model how vector representations of individual words are combined to generate vector representations of more complex sentences (e.g., Landauer & Dumais, 1997), and, more generally, seems to be the simplest and most parsimonious method of combining individual activation levels to determine the overall activation states involved in associative judgment. However, this technique ignores propositional structure. It is clear that the response *feminist and bank teller* is represented and evaluated differently to the response *feminist or bank teller* or the response *feminist and not bank teller* (again, see Carlson & Yates, 1989; Fisk, 2002; Nilsson et al., 2009), but the proposed model (unlike averaging rules, quantum judgment theory, and other accounts) cannot currently distinguish between these responses. Likewise although the use of cosine similarity to judge vector distances is standard in most applications of word vector models (again see Landauer & Dumais, 1997), it is symmetric, and thus cannot capture asymmetric aspects of similarity, or other asymmetries in judgment, such as those induced by order effects (and easily explained by quantum judgment models). However, it is important to note that representing complex propositional structure or permitting asymmetric similarity assessments is not outside the scope of the proposed approach, and in the discussion section of this paper I consider the possibility of incorporating the insights of existing theories into the proposed associative judgment model.

There are also effects related to the fallacies discussed in this article that models of associative judgment, including the one presented in this paper, cannot directly capture. One set of effects involves moderators of judgment fallacies. For example, prior work has found that presenting judgment problems in a frequency format reduces the incidence of the conjunction fallacy (Fiedler, 1988; Hertwig & Gigerenzer, 1999), and that base rate neglect is relatively infrequent when the base rates are implicitly learnt rather than described (Manis, Dovalina, Avis, & Cardoze, 1980; Medin & Edelson, 1988; see also Koehler, 1996 for a discussion). Additionally, there is much work on the effects of conversational norms on the interpretation of the conjunction term (e.g., Hertwig et al., 2008). It is useful to note that many other competing models of conjunction fallacies and base-rate neglect are also unable to account for these moderators: Overall, these moderators pertain to when different judgment strategies (such as those relying on associative processing vs. those relying on optimal probability calculations) are used, and explaining these moderators requires an analysis of the mechanisms involved in strategy selection (Marewski & Schooler, 2011; Payne, Bettman, & Johnson, 1988; Rieskamp & Otto, 2006).

Another set of effects involves judgment problems in which there is no description in the question and decision makers are asked to judge the probabilities of standalone events. Conjunction fallacies have been documented in these settings (see, e.g., problems in Tentori et al., 2013 and Gavanski & Roskos-Ewoldsen, 1991 not fit in the above analysis). This suggests that associative judgment may not be the only mechanism capable of generating the conjunction fallacy. Rather, some of the above mentioned theories, which focus on studying the rules involved in aggregating the probabilities of standalone events (rather than modeling the ways in which question descriptions activate response representations), may be better suited to explain this data.

However, ultimately the goal of this article is not to provide a single explanation for all observed judgment fallacies (indeed this would be impossible) but rather to examine the properties of an approach that is able to specify the associations involved across different types of judgment problems. Associative processing is an important feature of judgment, and there is value in formalizing associations so that theories of associative processing can be rigorously studied. However, it is not the only feature of judgment, and it is perfectly reasonable for this approach to not be able to describe every single finding related to the conjunction fallacy, as associations themselves are unable to explain every single finding related to the conjunction fallacy. Again, it is important to note that in trying to formalize associative judgment in this manner, the proposed approach has desirable properties that alternate explanations for the fallacies (such as those using configural weighting, quantum probability, confirmation, noisy recall, or probed recall) do not possess. Particularly, it is able to respond to natural-language questions without domain specific training or experimental testing, as it comes equipped with the representations necessary to apply it to most relevant judgment problems. Although the type of approach outlined in this article is common in other subfields in psychology, this is the first model in judgment and decision making research that has this property, and thus the first model to be able to make rigorous, a priori predictions about the judgments of decision makers, regardless of the specific (natural language) question that it is applied to. Overall, the analysis of the conjunction fallacy, and the data in Fischhoff and Bar-Hillel (1984), Gavanski and Roskos-Ewoldsen (1991), Shafir et al. (1990), and Tentori et al. (2013), as well as the analysis of our novel experimental data, does not only illustrate the predictive power of the vector space approach in high-level judgment, but also showcases its practical value and its methodological novelty.

Naturalistic Judgment Problems

This property of the proposed approach can be used to study more than the errors generated by associative judgment. The vector representations used in this article have vocabularies of between 300,000 and 3 million words, and are able to specify the associations between any questions and any responses that are composed of these words. Additionally, although the proposed approach is fundamentally based on associative relationships, the questions that we are able to answer need not be designed specifically to exploit associative judgment. Indeed, they need not even be designed by experimenters. The vector space-based approach can be applied to a wide range of naturalistic judgment problems, which have never before been used in psychological research.

Real-World Problems

I tested the ability of the proposed approach to describe judgments on three real-world question-answer data sets: geography quizzes obtained from the website www.about.com; elementary school multiple-choice problems used in the New York Regents examinations; and questions from the popular TV game show *Who Wants to be a Millionaire?* (WWTBAM).

These three data sets have a large number of questions, which cover a diverse array of topics. Each of the problems in these data sets is accompanied by a question and four possible response options, out of which one is correct and three are incorrect. We can use these questions and responses to test the proposed approach in a manner similar to the Linda problem described above. Particularly, each problem can be decomposed into five pieces of text: the question and the four responses. The associative strength between the words in the four responses and the words in the question, as assessed by our vector representations, can then be used to generate response predictions for the problem, so that the response with the strongest association with the words in the question is selected as the final prediction.

To test whether the vector space-based approach to uncovering associative relationships and making associative judgments is able to describe behavior in the above real-world question data sets, I performed three novel experiments (Studies 2–4). There were 100 participants recruited from MTurk in each of these experiments and these participants were given 30 questions, chosen at random from the geography quiz dataset in Study 2, the NY Regents dataset in Study 3, and the WWTBAM dataset in Study 4.

I first tested whether the proposed approach could predict participant accuracy on these questions. This was done with a linear regression. In this regression the proportion of participants selecting the correct response in each judgment problem was the dependent variable, and the normalized association between the correct response and the question was the independent variable. I performed this regression for the participant data in Studies 2–4, for each of the three vector space representations, and found that the normalized associations between the correct response and the question are positively associated with the proportion of participants answering the response correctly, in all nine cases ($p < .01$ for Word2Vec and GloVe and $p < .05$ for Eigenwords for the three experiments).

These results show that participants are more likely to be correct when the vector representations generate strong positive associations between the correct response and the question in consideration, implying that the vector representations are able to successfully predict participant accuracy across the different judgment problems. However, they do not show whether or not the specific response selected by the vector representations predicts the response chosen by the participants, independently of whether or not this response is correct. To test this, we can use a method similar to the one outlined in the previous section. Particularly, for each of the judgment problems in each of our data sets we can get the proportion of participants selecting each of the four feasible responses, and we can attempt to predict these response proportions using the normalized associations generated by the vector representations. Our test again involves a linear regression with random intercepts on the problem level. It also controls for the correctness of a particular response, and thus ensures that our results hold both

when the responses are correct and when these responses are incorrect. Performing this type of regression for our vector space representations again reveals that the associations generated by all vector representations have a positive relationship with the responses of participants ($p < .01$ for GloVe and Word2Vec for all three studies and for Eigenwords for Study 2 and Study 4, and $p < .05$ for Eigenwords for Study 3). A summary of these fits is provided in Table 2. Additional details regarding the question data sets and experimental methods are provided in the online supplemental materials.

Participant-Generated Problems

In this section I wish to test the predictive accuracy of the proposed approach on participant responses to questions generated by other participants from the same population. For this purpose I conducted a two-part experiment (Study 5). In the first part of the experiment I asked 50 participants on MTurk to generate five easy questions, five moderate questions, and five difficult questions each. Participants were also asked to list four candidate responses, with one response being correct and the remaining three being incorrect. In the second part, I asked another 100 MTurk participants to answer these judgment problems. Each participant in the second group answered 30 questions at random.

Can the proposed approach predict participant accuracy on these questions? For this I again use a linear regression with the proportion of participants selecting the correct response in each judgment problem as a dependent variable, and the predicted normalized association between the correct response and the question as the independent variable. This regression reveals that our three representations generate normalized associations that are positively correlated with the proportion of participants answering the response correctly ($p < .05$ for Word2Vec and Eigenwords and $p < .01$ for GloVe), indicating that the proposed approach quantitatively describes accuracy on participant-generated judgment problems.

As above we can expand upon this analysis by testing if the proposed approach also predicts the specific responses chosen by participants (independently of whether or not these responses are correct). Additionally, we can examine the effects of difficulty on these predictions, so as to determine whether the proposed approach's power is the same for easy questions, moderate questions, and hard questions. To test this I ran linear regressions with the

response proportions as the dependent variable, and the normalized associations as the independent variable. These regressions also included the correctness of the response ($correct = 1$ or 0 based on whether the response is correct), the difficulty of the question ($difficulty = 1, 2,$ or 3 based on whether the question is easy, moderate, or hard), and an interaction term between normalized associations and difficulty, as additional independent variables. This controls for correctness and difficulty, and the interaction term examines whether our ability to predict participant responses using normalized associations varies as a function of the difficulty of the questions. As above, the regressions also assumed random intercepts on the problem level.

Overall I found that the associations generated by our three vector representations have a positive relationship with the responses of participants ($p < .01$ for Word2Vec and GloVe, but, unlike in previous studies, $p > .10$ for Eigenwords). Additionally the predictive power of Word2Vec and GloVe decreases significantly as difficulty is increased, as evidenced by the negative interaction between the normalized association variable and the difficulty variable ($p < .05$ for Word2Vec and $p < .01$ for GloVe). This relationship is negative, but tiny and nonsignificant for Eigenwords. Once again, the results show that the proposed approach is able to quantitatively predict participant responses, though these predictions are less accurate for hard questions compared with easy questions, and less accurate for the Eigenwords representations. A summary of fits is provided in Table 3. Additional details regarding the experimental methods are provided in the online supplemental materials.

Comparison With Recognition

One limitation of the above tests is that they do not compare the proposed approach against a competing theory of judgment. One reason for this is that there is no other existing psychological theory, approach, or technique that is able to make a priori predictions for the types of natural language problems studied in this article. One exception to this is the recognition heuristic, which involves the use of recognition memory to predict participant responses (Gigerenzer & Goldstein, 1996; Goldstein & Gigerenzer, 2002; Pachur & Hertwig, 2006; Schooler & Hertwig, 2005). Perhaps the best known illustration of this is the city-size judgment task (Gigerenzer & Goldstein, 1996). In this task decision makers are required to judge which of two German cities is

Table 2
Summary of Model Fits of the Three Vector Representations to the Three Real-World Question Datasets

Data	Word2Vec			Eigenwords			GloVe		
	β	z	R^2	β	z	R^2	β	z	R^2
Study 2 (Geo Quiz)	1.73	5.77**	.06	.34	2.26**	.04	.67	3.83**	.05
Study 3 (NY Regents)	1.01	4.42**	.86	.27	2.18*	.84	.34	2.70**	.85
Study 4 (WWTBAM)	1.21	4.97**	.42	.37	2.88**	.40	.50	3.99**	.42
Average R^2			.45			.43			.44

Note. The fits involve random-effects linear regressions to predict the responses of participants using the normalized associations generated by the vector representations. These regressions also control for whether or not the response in consideration is correct or not.

* $p < .05$. ** $p < .01$.

Table 3
Summary of Model Fits of the Three Vector Representations to the Participant Generated Questions in Study 5

Representations	β	z	β_{int}	z_{int}	R^2
Word2Vec	1.91	2.92**	-.58	-2.04*	.76
Eigenwords	.19	.40	-.01	-.04	.66
GloVe	2.33	5.03**	-.80	-3.83**	.77

Note. The fits involve random-effects linear regressions, controlling for question difficulty and response correctness. They also include an interaction between difficulty and normalized associations. Here β and z correspond to the main effect of normalized associations, and β_{int} and z_{int} correspond to the interaction between normalized associations and difficulty.

* $p < .05$. ** $p < .01$.

the largest. Gigerenzer and Goldstein suggested that decision makers use recognition to solve this problem so that if they are able to recognize one of the two cities, but not the other, they infer that the recognized city is the largest. An extension of the recognition heuristic in which fluency, a variable that captures the overall ease with which the recognition judgment is made, has also been examined (Hertwig, Herzog, Schooler, & Reimer, 2008; Marewski & Schooler, 2011; Schooler & Hertwig, 2005). More important, for our purposes, both the recognition and the fluency of an object (e.g., a city) can be specified by frequency of mention in natural language (Goldstein & Gigerenzer, 2002).

If the responses to a question are simple words and phrases that refer to a single object or concept, such as a city, it could be possible to use the recognition or fluency of that object (captured through its frequency of occurrence in a relevant natural language dataset) to make a prediction regarding the likelihood of decision makers choosing that object as their response. This would not only facilitate a comparison between the vector space approach and a competing approach, but would also ensure that the predictive power of the proposed approach does not stem from a correlation between association and recognition, that is, that the results outlined above are not confounded by the recognizability of the responses in consideration.

I attempt such a comparison using judgment problems involving cities. Cities are, as discussed above, the objects that the recognition heuristic has most commonly been applied to. Indeed, Goldstein and Gigerenzer (2002) and Marewski and Schooler (2011) have already suggested ways of making a priori predictions regarding the recognizability or fluency of these cities, using metrics such as frequency of mention in news media and search frequency on online search engines (see also Anderson & Schooler, 1991). For any given judgment problem whose responses are a list of cities, we can use this metric to quantify the recognizability of each city, and compare this with the association between that city and the content of the question.

To perform this test, I again ran a two-part experiment (Study 6). In this experiment I asked one MTurk participant to generate one fact each for 16 large U.S. cities. In the second part of our experiment I used these 16 facts about the cities to generate 48 multiple-choice judgment problems, with four options each. Each of these judgment problems offered a city fact as a description and then asked the respondents to select a city that matched the description offered. The four options included the city to which the

fact applied as well as three other cities randomly chosen from the list of 16 cities. These 48 judgment problems were then offered to 150 other MTurk participants, who were given 10 randomly chosen judgment problems each.

We can use the responses generated in this study to compare our vector space approach with the recognition heuristic. For the purposes of this Test I formalized recognizability using Google Trends search frequency (as in Marewski & Schooler, 2011) and city occurrence frequency in the Google books corpus (as in Goldstein & Gigerenzer, 2002). These variables, which I label *search frequency* and *n-gram frequency*, provide us with two similar but nonidentical ways of formalizing the recognizability of each city. As with association, city recognizability is normalized for each question with a zero-parameter soft-max transform.

Once again tests can be done using simple linear regressions in which the proportion of participants selecting each of the cities as answers to the 48 questions serves as the dependent variable. The independent variables in this case are the normalized association of each city with the question at hand as well as the normalized recognizability of each city. Additionally, I permit random effects on the question level. The first set of regressions that I consider involve running both of the independent variables simultaneously so as to test whether the positive relationship between association and participant responses persists once the recognizability of these responses is controlled for. As there are three ways of formalizing association and two ways of formalizing recognizability this leads to a total of six regressions. Performing all six regressions I find that associations are positively related to participant responses ($p < .01$ for each of our three vector representations for each regression), and that recognizability is either not statistically related to these responses ($p > .5$ for our two recognizability variables for five of the regressions) or negatively related to these responses ($p < .01$ when search frequency is regressed alongside GloVe's representations). This indicates that the recognizability of a response is not a confound when testing the descriptive power of the proposed approach.

The second set of regressions that we can run involves testing the independent variables separately, so as to examine their individual ability to predict participant responses. Here we have five separate regressions, three for association and two for recognition. Again I find that all three of the vector representations generate associations that are positively related to participant responses ($p < .01$ for all three representations). In contrast, the two ways of specifying recognition do not show any significant relationship with the responses of participants ($p > .90$ for both recognizability variables). These regressions are summarized in Table 4, which also shows the R^2 values generated by each of these five regressions. As can be seen in this table, the Word2Vec and GloVe representations are able to capture a very large proportion of the variance in participant responses, indicating that these methods of specifying associations are particularly good at predicting participant choices. Additional details regarding the experimental methods are provided in the online supplemental materials.

Note that these tests do not necessarily imply that the recognition heuristic is not a good account of judgment. The proponents of this heuristic have explicitly stated that recognition is likely to be used only when it correlates with the criterion variable (Gigerenzer & Goldstein, 1996; Goldstein & Gigerenzer, 2002; Marewski & Schooler, 2011; Schooler & Hertwig, 2005). Thus, recog-

Table 4
Summary of Model Fits to Study 6

Model	β	z	R^2
Word2Vec	13.49	13.31*	.48
Eigenwords	7.21	2.84*	.04
GloVe	16.62	19.22*	.66
Search Frequency	-.01	-.08	.00
N-gram Frequency	-11.39	-.08	.00

Note. The fits involve random-effects linear regressions to predict the responses of participants using either normalized associations generated by each of the three vector representations or the normalized recognizability generated by the search frequency and n-gram frequency variables. These regressions are all performed separately.

* $p < .01$.

inition is a useful cue for predicting city size and is used when participants are asked to judge which of two cities is the largest. It is not necessarily a useful cue for assessing other facts about cities, and thus may not be particularly active in the judgment problems I examine above. The use of recognition is merely to permit a comparison between the proposed approach and an existing approach, and the recognition heuristic is the only existing approach, to my knowledge, that is able to make precise a priori predictions for the types of natural language judgment problems studied in this article.

Event Probabilities

Thus far I have examined multiple-choice problems in which each question is accompanied by a set of feasible responses. There is also, however, a related type of task in which associations are active. This involves the assessment of real-world event probabilities, without a specific multiple choice format. For example, people are often asked to attach explicit probabilities to various outcomes in current affairs or popular culture. These types of events are often of the form “X happens to Y” (e.g., *an earthquake occurs in Japan*), and theories of associative judgment predict that the association between X and Y (e.g., *earthquake* and *Japan*) is used by individuals to judge the probabilities of these types of events.

We can modify the proposed approach to predict these probability judgments, by using the vector space representations to specify the associations between the various components of the events. Particularly, for events of the form “X happens to Y” these probabilities can be predicted by examining the proximity of the vectors corresponding to X and Y. In Studies 7 and 8 I tested this idea with probability judgments about different countries and different famous people respectively. There were 200 MTurk participants each in Studies 7a–7d and 100 MTurk participants each in Studies 8a–8d. For each of the countries offered to the participants, they were asked to assess the probability that the country would experience a terrorist attack in the next week (Study 7a), be in a state of war at the start of 2016 (Study 7b), experience an earthquake over the next year (Study 7c), or experience an epidemic over the next year (Study 7d). Each participant was given a list of 30 countries chosen at random from the 193 countries that are members of the United Nations. Likewise, for each of the people offered to the participants, they were asked to assess the

probability that the person would become the U.S. president in 2020 (Study 8a), win a Nobel Prize in 2020 (Study 8b), win a Grammy Award in 2020 (Study 8c), or win an Academy Award in 2020 (Study 8d). Each participant in Studies 8a–8d made judgments about 30 people chosen at random from our list of 50 people.

The average probabilities obtained in Studies 7 and 8 were predicted using only the Word2Vec representations. This is because of the large vocabulary of these representations, which contains not only individual words but also combinations of words, including country and person names like *United States* and *Jon Stewart*. In contrast, the vocabularies of our Eigenwords and Glove representations only have single words. For Studies 7a–7d, I used the cosine similarity between *terrorism*, *war*, *earthquake*, and *epidemic* and the words corresponding to the 193 countries to predict the probabilities that participants assign to the disasters happening in the countries. Likewise I used the cosine similarity between *President*, *Nobel Prize*, *Grammy Award*, and *Academy Award* and the names of the 50 famous people to predict the probabilities that participants assign to the people winning these awards in Studies 8a–8d.

For each of the events across the 8 studies we have both the average probability assigned to the event by the participants, and the association between the words in the event, specified by the Word2Vec representations. A first step in the analysis is examining the correlation between associations and the average estimates of participants. Pearson’s correlation reveals positive correlations between these two variables in each of our studies ($p < .01$). Scatter plots displaying the relationship between associations and participant judgments can be seen in Figure 4a–h, and the correlations outlined here are summarized in Table 5. For a more rigorous test, I again considered a linear regression. This transforms the cosine similarity measure of association, which ranges from -1 to 1 , into a probability judgment scale, which ranges from 0 to 100 . After performing this regression for the eight studies I unsurprisingly found highly significant relationships between associations and participant judgments ($p < .01$ for each of the studies). More details about these fits are provided in Table 5. Additionally, details regarding the experimental methods are provided in the online supplemental materials.

A Note on Accuracy

The results in the above sections illustrate the ability of the proposed approach to describe judgments in real-world problems. These problems are not only useful for studying the external validity of theories of associative judgment, but also in examining the adaptive value of this type of judgment. Typically, questions such as the Linda problem are used to highlight the tendency for decision makers to make errors (Gilovich et al., 2002; Kahneman & Tversky, 1973; Tversky & Kahneman, 1983). However, it is clear that judgment processes need to be, on aggregate, beneficial, as they are unlikely to be adopted if they always lead to incorrect responses (Anderson, 1990; Gigerenzer & Todd, 1999; Oaksford & Chater, 2007; Simon, 1990).

There is some evidence that vector space semantic models are adaptive. For example, Landauer and Dumais (1997) showed that latent semantic analysis, a prominent vector space model, is able to answer judgment problems in English language tests accurately

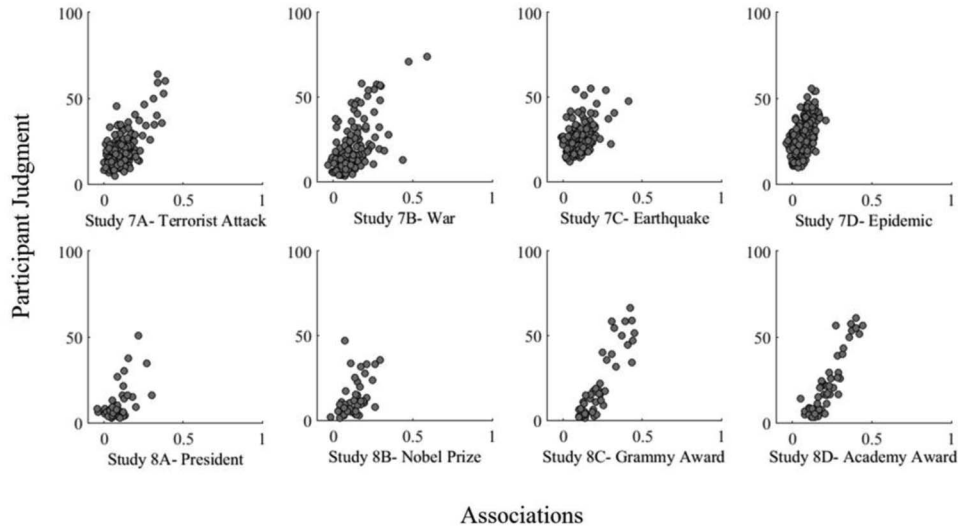


Figure 4. (a–h) Scatter plots of the word associations (in terms of cosine similarity) generated by the Word2Vec representations and average participant probability estimates (in terms of percentage) for the events in Studies 7a–7d and 8a–8d.

(also see Griffiths et al., 2007; Jones & Mewhort, 2007 for similar results; as well as Bullinaria & Levy, 2007; Jones et al., 2015 for a review). Although these tests are primarily linguistic in nature, that is, they involve judgments of word meaning, similarity, and so on, the overall accuracy of vector space representations on these tests supports the idea that the use of these representations is beneficial for decision making. It also suggests that these representations may also be able to accurately answer the types of high-level judgment problems studied in this paper, despite the fact that these judgment problems involve much more complex (often nonlinguistic) inferences.

So, how well do the associations generated by the vector representations answer judgment problems? We can test this by applying them to the geography quiz, New York Regents' and Who Wants to be a Millionaire question data sets for Studies 2–4, and using them not to predict participant responses to questions in these data sets but the correct responses in these data sets. Doing so shows that all three of our representations significantly ($p < .01$ for Word2Vec and GloVe, and $p < .05$ for Eigenwords) outper-

form random choice (25% accuracy) by correctly answering up to 42% of the questions in these data sets. This is shown in Figure 5. We can also test the accuracy of the proposed approach on our participant-generated problems obtained in Study 5. For these problems, I find that the Word2Vec and GloVe vector representations outperform random choice for all difficulty levels ($p < .01$), though Eigenwords only does so for easy questions. These accuracy rates are shown in Figure 6.

The adaptive value of the proposed approach can also be tested using the city size task. As outlined above, decision makers in this task are required to judge which of two German cities is the largest. Gigerenzer and Goldstein (1996) suggested that decision makers use recognition to solve this problem and show that for a list of 84 large German cities, a heuristic that relies only on recognition can achieve an accuracy rate of up to 65% (a fluency-based extension of the recognition heuristic can lead to further improvements in accuracy, as shown in Marewski & Schooler, 2011). Of course recognition is not the only judgment process capable of answering the above question successfully. It may be the case that associative processing plays an important role in facilitating high accuracy in the city-size judgment task and related problems. Indeed, the proposed approach is able to answer questions like *Which of the following has the biggest population?* Response 1: *Hamburg*; Response 2: *Cologne*, by examining the strength of associations between the question and the two responses.

After applying the proposed approach to the above question and to the cities in Gigerenzer and Goldstein I found that the GloVe and Word2Vec vector representations are able to achieve accuracy rates of 72% and 64% in paired judgment, far outperforming the random chance rate of 50%. This is comparable with the accuracy rates achieved using only recognition. In contrast to this, Eigenwords representations only achieve accuracy rates of 54%, suggesting that not all ways of building vector representations are useful for this task. As a second Test I also ran a linear regression

Table 5

Summary of Linear Model Fits of the Associations Generated by Word2Vec Representations and the Aggregate Judgments of Participants in Studies 7a–d and Studies 8a–d

Study	Event	Target	ρ	β	t	R^2
7a	Terrorist attack next week	Country	.62	84.48	10.82	.38
7b	Being at war at start of 2016	Country	.63	94.29	11.29	.40
7b	Earthquake next year	Country	.44	53.60	6.72	.19
7d	Epidemic next year	Country	.57	118.26	9.46	.32
8a	US president in 2020	Person	.59	85.56	17.30	.35
8b	Nobel Prize in 2020	Person	.47	75.48	20.40	.23
8c	Grammy Award in 2020	Person	.89	161.74	13.28	.79
8d	Academy Award in 2020	Person	.90	163.82	14.12	.81

Note. All coefficients are significant at $p < .01$.

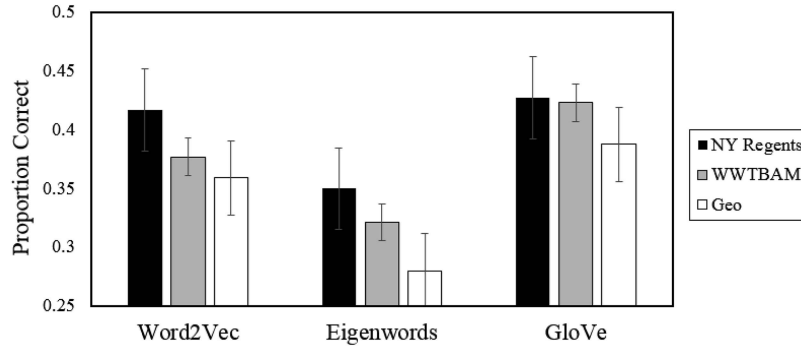


Figure 5. The proportion of questions in the three real-world question-answer data sets for which the vector representations specify the correct answer. Note that 25% is the accuracy expected from a model that makes random selections.

with the ordinal rank of the size of the city as the dependent variable and the associations generated by the vector representations as the independent variables. Using this method I found that all three of the vector representations generated associations between the name of the city and the words in the judgment question, that are statistically related to the rank of the city in terms of city size ($p < .01$ for all three approaches). These results, summarized in Table 6, demonstrate that strength of association is a cue that is positively correlated with accuracy, and should be taken into consideration by decision makers.

General Discussion

In this article I have applied research on semantic memory and computational linguistics to judgment and decision making to test a novel approach to modeling associative judgment. The proposed approach uses the semantic relatedness specified by vector space semantic models (particularly those in Dhillon et al., 2011; Mikolov et al., 2013; Pennington et al., 2014) to quantify the associations involved in judgment (Kahneman, 2003; Kahneman & Frederick, 2002; Sloman, 1996), and in turn to predict the responses of participants in a large range of existing and novel, experimental and real-world, and experimenter-generated and participant-generated judgment problems. I have found that vector

space representations are able to successfully describe both qualitative and quantitative patterns underlying participant responses, across all these different types of problems, and that they are able to provide a good account of the associative processes at play in judgment.

Formalizing Associative Judgment Processes

Theories of associative judgment were initially proposed to explain biases such as the conjunction fallacy (Kahneman, 2003; Kahneman & Frederick, 2002; Sloman, 1996; Tversky & Kahneman, 1983). In the above sections, I have shown how the proposed approach to formalizing associative judgment, is able to account for a large number of findings involving this fallacy, as well as its various moderators, such as those pertaining to the role of typicality, averaging, and confirmation (e.g., Shafir et al., 1990; Gavanski & Roskos-Ewoldsen, 1991; Tentori et al., 2013). Additionally, these sections have applied this approach to related fallacies such as base rate neglect (Fischhoff & Bar-Hillel, 1984). More important, these explanations are not only qualitative but also quantitative, so that the proposed approach is able to make successful a priori predictions regarding the precise probability assignments and response proportions observed in prior experimen-

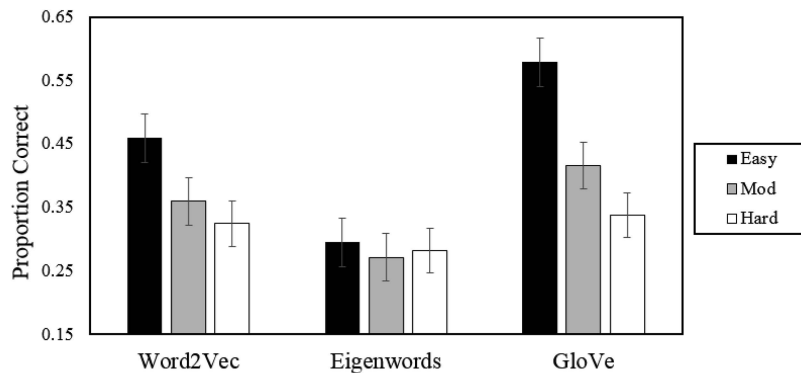


Figure 6. The proportion of easy, moderate, and hard questions in participant-generated problems (Study 5) for which the vector representations specify the correct answer. Note that 25% is the accuracy expected from a model that makes random selections.

Table 6
Summary of Linear Regressions of City Size Rank (With Lower Ranks Indicating Bigger Cities) on Normalized Associations Generated by the Three Sets of Vector Representations

Representations	β	z	R^2
Word2Vec	-139.68	-3.60	.34
Eigenwords	-.29	-3.64	.15
GloVe	-128.45	-6.63	.36

Note. All coefficients are significant at $p < .01$.

tal data sets involving these fallacies, as well as in novel participant-generated fallacy problems.

Despite its unique descriptive power, the vector space approach is not without limitations. As discussed above, this approach is not able to explain a number of findings regarding the moderators of these biases, such as why presenting judgment problems in a frequency format reduces the incidence of the conjunction fallacy, why base rate neglect is relatively infrequent when the base rates are implicitly learnt rather than described, or how conversational norms affect the interpretation of the conjunction term (Fiedler, 1988; Hertwig et al., 2008; Hertwig & Gigerenzer, 1999; Koehler, 1996). This indicates that a complete theory of the conjunction fallacy and related biases needs to incorporate other mechanisms besides the ones studied in this paper. That said, by formalizing associative judgment, one of the most prominent explanations of these biases, and, in turn, specifying a way to apply this theory's predictions across a wide range of settings, the proposed approach adds greatly to our understanding of the errors and cognitive illusions involved in judgment. In doing so, it also answers the call for more rigorous theories of associative judgment (Gigerenzer, 1996, 1998; Gigerenzer & Regier, 1996).

Naturalistic Judgment

To test the external validity of the vector space approach to modeling judgment, this article examined participant-generated judgment problems, problems obtained from real-world quiz and exam compilations, and event-probability judgment problems. These questions span a countless range of topics and presentation formats, and closely resemble the types of naturalistic judgment problems faced by participants outside of the laboratory. In all of these tasks, this article found that the proposed approach is able to successfully predict participant responses. Now, there are some limitations to using the proposed approach for modeling naturalistic judgments, as people use much more than associations when taking exams or when judging events. However, by showing that associations are strongly correlated with participant responses, this paper provides one of the strongest possible tests of associative judgment, and of the vector space-based approach to modeling associative judgment. This approach does not only provide a good description of behavior in simplistic artificial tasks where responses are explicitly chosen to be associated with the question, but also in the more naturalistic tasks people encounter in real-world decision making settings.

Adaptive Judgment

A related contribution of this paper is in showing that associative judgment (as specified by the proposed vector space approach)

does not only generate error. Rather it is capable of making real-world judgments with an above chance level of accuracy. As processing associations is relatively automatic and effortless, these results imply that association-based judgment can be seen as being adaptively rational, and not only a source of bias. Ultimately, the associative strength between a question and a response is a good cue for assessing the correctness of the response, in the real world, and should thus be used by decision makers.

The importance of adaptive judgment has been stressed in the works of Gigerenzer and colleagues, whose research program has explicitly focused on understanding the processes underlying accurate behavior and on studying how judgment relies on memory representations that reflect the structure of the environment (Gigerenzer & Todd, 1999; Gigerenzer & Gaissmaier, 2011). Although similar insights have been made with regards to linguistic judgment (e.g., Landauer & Dumais, 1997), these claims are seldom extended to high-level associative judgment. For this reason, the biased (associative) processes specified by Kahneman, Tversky, and coauthors, and the adaptive processes specified by Gigerenzer and coauthors are often studied separately. However, it is clear that both perspectives warrant merit and that a complete account of judgment needs to be able to capture both error and intelligence. This article provides some steps toward such an account. In doing so, it highlights the importance of an integrative approach to studying judgment, one that combines insights from multiple research traditions, to provide a richer account of heuristic decision making and its relationship with rationality.

Representation and Judgment

One of the reasons that the proposed approach is able to perform the above tests is because it explicitly represents the type of information that associative judgment utilizes: that is, it knows about both a wide range of concepts, and the associative relationships between these concepts. This aspect of information representation is a relatively neglected area in judgment and decision making research. Although this field has a number of powerful theories regarding how information is used in judgment (Busemeyer et al., 2011; Dougherty et al., 1999; Gigerenzer & Gaissmaier, 2011; Gilovich et al., 2002; Hammond & Stewart, 2001; Kahneman & Tversky, 1973; Reyna et al., 2003; Shah & Oppenheimer, 2008; Tversky & Koehler, 1994; see also Weber & Johnson, 2009 and Hastie, 2001), and has also recently made important strides in understanding the processes involved in learning this information (e.g., Broder & Gaissmaier, 2007; Dougherty et al., 1999; Juslin & Persson, 2002; Lagnado et al., 2006), work in this field has seldom specified, a priori, what this information actually is.

Consider, for example, the Linda problem (Tversky & Kahneman, 1983). Responses to this problem are often attributed to the computations performed by associative processes. However, proponents of this theory do not attempt to explain how decision makers associate a woman of Linda's description with a feminist. The connection between Linda and a feminist is first grounded in the experimenter's intuition, and occasionally established by asking participants to make similarity ratings between Linda's description and that of a feminist. It does not stem from the judgment theory itself.

Fortunately, the problem of representation has been tackled in other subfields of psychology and cognitive science. This work proposes that people's representation of words depends on the statistical structure of the environment in which these words occur (Firth, 1957 and Harris, 1954). Studying the distribution of words in the types of settings people encounter on a day-to-day basis can uncover the representations that people have of everyday words, and in turn the associations between these words, and the objects and concepts they represent. Models that build representations using the distribution of words often characterize each word in their vocabulary as a vector in a multidimensional space, with the proximity between two vectors corresponding to the semantic relatedness or association between their words (Dhillon et al., 2011; Griffiths et al., 2007; Jones & Mewhort, 2007; Kwantes, 2005; Landauer & Dumais, 1997; Lund & Burgess, 1996; Mikolov et al., 2013; Pennington et al., 2014). These models are typically trained on very large natural language text corpora, and subsequently have large vocabularies, which can be used to specify representations for a large number of words.

The insights of vector space models have mostly been used to understand semantic memory and predict effects pertaining to similarity judgment, priming, recall, and related psycholinguistic phenomena (Bullinaria & Levy, 2007; Jones et al., 2015). The approach proposed in this article is motivated by the success of this work, and adopts three state-of-the-art techniques for building vector representations (Dhillon et al., 2011; Mikolov et al., 2013; Pennington et al., 2014) to the study of judgment. Its results suggest that the key mechanisms used to build representations in semantic memory research can also describe the representations underlying associative judgment. Decision makers believe Linda is a feminist because instances of feminism frequently co-occur with intellectual inclination and a desire for social justice, in the descriptions of various individuals. By training vector space models on natural language data sets with these descriptions, it is possible to learn representations that capture these associations, and subsequently predict, a priori, the probability attached to Linda being a feminist.

Extensions

This article has only considered three techniques for building vector space representations, and additionally has used prebuilt vector representations generated by these techniques, rather training its own representations. This has been primarily because of computational concerns. To successfully apply the proposed approach to predict judgments across different types of problems, it is necessary for its vector representations to have very large vocabularies. The goal of this article is to describe judgment using the broader vector space-based approach, rather than to propose or test the power of a single new or existing vector space model. Thus, the use of these existing representations has not detracted from the claims made in this paper. That said, it would be useful, in future work, to attempt to apply these techniques to a standardized natural language corpus, to provide a rigorous test of their relative power. Such a test could also build representations using other existing computational techniques, such as those involving word-based semantic spaces (Lund & Burgess, 1996), singular value decomposition (Landauer & Dumais, 1997), convolution and superposition mechanisms (Jones & Mewhort, 2007), and topic

models (Griffiths et al., 2007), many of which provide a good account of the representations at play in semantic memory tasks (see Bullinaria & Levy, 2007; Jones et al., 2015).

A related extension would involve training the semantic models on individual-level natural language corpora. This would allow them to make individual-level predictions and describe individual-level differences, which is not currently possible using the proposed approach. Such an extension would shed light on a number of issues of importance in judgment research, including how different knowledge representations lead to different judgment biases, why experts outperform nonexperts in some settings but not in others, and the ways in which we could improve judgment by modifying individual knowledge representations. It would also allow for the examination of cultural, linguistic, gender, and age-related differences in judgment with more rigor. Training individual-level representations may not be immediately feasible; however, the increased digitization of information may make this a possibility in the near future. In the meantime, it is possible to examine some individual differences in judgment by using vector space models trained on the types natural language data that different demographics are more or less likely to be exposed to.

There are other ways to improve upon the results of this article. These do not involve the use of different representations, but rather different manipulations of the representations to make judgments. For example, the strength of association between a question and a response used in this article, cosine similarity, is symmetric. However, this type of symmetry is often not observed in similarity-based judgment (Tversky, 1977). Likewise, by using a bag-of-words approach, the model ignores word order. Word order, however, plays a key role in a number of different types of judgments, most notably in judgments regarding concept relations and analogy. Word order is also necessary to learn, represent, and respond to most types of logic structure. Although violations of the conjunction fallacy, as with the Linda problem, indicate that decision makers may not always process logical structure, it is nonetheless necessary for a judgment model to perceive and utilize basic logical operators, such as disjunction and negation. After all, the response *feminist and bank teller* is represented and evaluated differently to the response *feminist or bank teller* (Carlson & Yates, 1989; Fisk, 2002; Nilsson et al., 2009).

One promising approach to solving this problem involves quantum judgment (Busemeyer et al., 2011; Trueblood & Busemeyer, 2011), which relies on a closely related vector-based representation format. Unlike the work in this article, theories of quantum judgment are not concerned with what these representations are, but rather how manipulations of these vector representations, based on quantum information processing principles, could be used to account for the effects of word order as well as different logical connectives (e.g., conjunctions vs. disjunctions) in judgment problems. It would be possible to combine the word vectors used in the proposed approach, with the vector projection-based conjunction, disjunction, and negation rules, specified by theories of quantum judgment, as suggested already by Kintsch (2014); see also Aerts & Czachor, 2004; Kintsch, 2001; Mitchell & Lapata, 2010). However, a limitation of assuming this type of method for building vectors is that it may not be easily applicable to settings in which the questions and responses involve more than just pairs of words connected through logical operators.

Conclusion

There are a number of desirable theoretical approaches to modeling judgment, yet these approaches are limited to studying only the processes involved in judgment. They do not specify the representations that judgment processes utilize. This article has shown how vector space models developed in semantic memory research and computational linguistics can be used to specify some of these representations, and rigorously model associative judgment. This integrative approach has allowed us to make successful a priori qualitative and quantitative predictions for a large variety of judgment problems, including problems used in existing research, new participant-generated problems, problems obtained from real-world question compilations, and real-world event probability judgment problems.

Although vector space semantic models are commonly used to predict responses in linguistic and semantic memory tasks, their ability to do so with regards to associative judgment suggests a new way of studying high-level judgment processes. By possessing both the same (association-based) cognitive mechanisms as humans and the same (vector-based) information representations as humans, the proposed approach has the ability to respond in a human-like manner to a very large array of judgment problems. In testing and verifying the power of this approach, this article contributes to a heightened degree of formalism in decision modeling, and illustrates how insights from semantic memory research and related areas can be used to build a new class of powerful, flexible, domain-general theories of judgment and decision making.

References

- Aerts, D., & Czachor, M. (2004). Quantum aspects of semantic analysis and symbolic artificial intelligence. *Journal of Physics A, Mathematical and General*, 37, L123–L132. <http://dx.doi.org/10.1088/0305-4470/37/12/L01>
- Anderson, J. R. (Ed.). (1990). *The adaptive character of thought*. Hove, England: Psychology Press.
- Anderson, J. R., & Schooler, L. J. (1991). Reflections of the environment in memory. *Psychological Science*, 2, 396–408. <http://dx.doi.org/10.1111/j.1467-9280.1991.tb00174.x>
- Barbey, A. K., & Sloman, S. A. (2007). Base-rate respect: From ecological rationality to dual processes. *Behavioral and Brain Sciences*, 30(03), 241–254.
- Bar-Hillel, M., & Neter, E. (1993). How alike is it versus how likely is it: A disjunction fallacy in probability judgments. *Journal of Personality and Social Psychology*, 65, 1119–1131. <http://dx.doi.org/10.1037/0022-3514.65.6.1119>
- Bröder, A., & Gaissmaier, W. (2007). Sequential processing of cues in memory-based multiattribute decisions. *Psychonomic Bulletin & Review*, 14, 895–900. <http://dx.doi.org/10.3758/BF03194118>
- Bullinaria, J. A., & Levy, J. P. (2007). Extracting semantic representations from word co-occurrence statistics: A computational study. *Behavior Research Methods*, 39, 510–526. <http://dx.doi.org/10.3758/BF03193020>
- Busemeyer, J. R., Pothos, E. M., Franco, R., & Trueblood, J. S. (2011). A quantum theoretical explanation for probability judgment errors. *Psychological Review*, 118, 193–218. <http://dx.doi.org/10.1037/a0022542>
- Carlson, B. W., & Yates, J. F. (1989). Disjunction errors in qualitative likelihood judgment. *Organizational Behavior and Human Decision Processes*, 44, 368–379. [http://dx.doi.org/10.1016/0749-5978\(89\)90014-9](http://dx.doi.org/10.1016/0749-5978(89)90014-9)
- Costello, F. J. (2009). How probability theory explains the conjunction fallacy. *Journal of Behavioral Decision Making*, 22, 213–234. <http://dx.doi.org/10.1002/bdm.618>
- Costello, F., & Watts, P. (2014). Surprisingly rational: Probability theory plus noise explains biases in judgment. *Psychological Review*, 121, 463–480. <http://dx.doi.org/10.1037/a0037010>
- Crupi, V., Fitelson, B., & Tentori, K. (2008). Probability, confirmation, and the conjunction fallacy. *Thinking & Reasoning*, 14, 182–199. <http://dx.doi.org/10.1080/13546780701643406>
- Crupi, V., & Tentori, K. (2016). Noisy probability judgment, the conjunction fallacy, and rationality: Comment on Costello and Watts (2014). *Psychological Review*, 123, 97–102. <http://dx.doi.org/10.1037/a0039539>
- Dhillon, P., Foster, D. P., & Ungar, L. H. (2011). Multi-view learning of word embeddings via cca. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems* (pp. 199–207). New York, NY: Curran Associates, Inc.
- Dhillon, P., Foster, D., & Ungar, L. (2015). Eigenwords: Spectral word embeddings. *Journal of Machine Learning Research*, 16, 3035–3078.
- Dhillon, P., Lu, Y., Foster, D. P., & Ungar, L. (2013). New subsampling algorithms for fast least squares regression. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Q. Weinberger (Ed.), *Advances in neural information processing systems* (pp. 360–368). New York, NY: Curran Associates, Inc.
- Dougherty, M. R., Gettys, C. F., & Ogden, E. E. (1999). MINERVA-DM: A memory processes model for judgments of likelihood. *Psychological Review*, 106, 180–209. <http://dx.doi.org/10.1037/0033-295X.106.1.180>
- Erev, I., Wallsten, T. S., & Budescu, D. V. (1994). Simultaneous over- and underconfidence: The role of error in judgment processes. *Psychological Review*, 101, 519–527. <http://dx.doi.org/10.1037/0033-295X.101.3.519>
- Evans, J. S. B. (2008). Dual-processing accounts of reasoning, judgment, and social cognition. *Annual Review of Psychology*, 59, 255–278.
- Fantino, E., Kulik, J., Stolarz-Fantino, S., & Wright, W. (1997). The conjunction fallacy: A test of averaging hypotheses. *Psychonomic Bulletin & Review*, 4, 96–101. <http://dx.doi.org/10.3758/BF03210779>
- Fiedler, K. (1988). The dependence of the conjunction fallacy on subtle linguistic factors. *Psychological Research*, 50, 123–129. <http://dx.doi.org/10.1007/BF00309212>
- Firth, J. R. (1957). *Papers in Linguistics*. London, England: Oxford University Press.
- Fischhoff, B., & Bar-Hillel, M. (1984). Diagnosticity and the base-rate effect. *Memory & Cognition*, 12, 402–410. <http://dx.doi.org/10.3758/BF03198301>
- Fisk, J. E. (2002). Judgments under uncertainty: Representativeness or potential surprise? *British Journal of Psychology*, 93, 431–449. <http://dx.doi.org/10.1348/000712602761381330>
- Franco, R. (2009). The conjunction fallacy and interference effects. *Journal of Mathematical Psychology*, 53, 415–422. <http://dx.doi.org/10.1016/j.jmp.2009.02.002>
- Gavanski, I., & Roskos-Ewoldsen, D. R. (1991). Representativeness and conjoint probability. *Journal of Personality and Social Psychology*, 61, 181–194. <http://dx.doi.org/10.1037/0022-3514.61.2.181>
- Gigerenzer, G. (1996). On narrow norms and vague heuristics: A reply to Kahneman and Tversky (1996). *Psychological Review*, 103, 592–596. <http://dx.doi.org/10.1037/0033-295X.103.3.592>
- Gigerenzer, G. (1998). Surrogates for theories. *Theory & Psychology*, 8, 195–204. <http://dx.doi.org/10.1177/0959354398082006>
- Gigerenzer, G., & Gaissmaier, W. (2011). Heuristic decision making. *Annual review of psychology*, 62, 451–482.
- Gigerenzer, G., & Goldstein, D. G. (1996). Reasoning the fast and frugal way: Models of bounded rationality. *Psychological Review*, 103, 650–669. <http://dx.doi.org/10.1037/0033-295X.103.4.650>

- Gigerenzer, G., & Regier, T. (1996). How do we tell an association from a rule? Comment on Sloman (1996). *Psychological Bulletin*, *119*, 23–26. <http://dx.doi.org/10.1037/0033-2909.119.1.23>
- Gigerenzer, G., & Todd, P. M. (1999). *Simple heuristics that make us smart*. New York, NY: Oxford University Press.
- Gilovich, T., Griffin, D., & Kahneman, D. (Eds.). (2002). *Heuristics and biases: The psychology of intuitive judgment*. New York, NY: Cambridge University Press. <http://dx.doi.org/10.1017/CBO9780511808098>
- Goldstein, D. G., & Gigerenzer, G. (2002). Models of ecological rationality: The recognition heuristic. *Psychological Review*, *109*, 75–90. <http://dx.doi.org/10.1037/0033-295X.109.1.75>
- Griffiths, T. L., Steyvers, M., & Tenenbaum, J. B. (2007). Topics in semantic representation. *Psychological Review*, *114*, 211–244. <http://dx.doi.org/10.1037/0033-295X.114.2.211>
- Hammond, K. R., & Stewart, T. R. (2001). *The essential Brunswick: Beginnings, explications, applications*. New York, NY: Oxford University Press.
- Harris, Z. S. (1954). Distributional structure. *Word*, *2*, 146–62.
- Hastie, R. (2001). Problems for judgment and decision making. *Annual Review of Psychology*, *52*, 653–683. <http://dx.doi.org/10.1146/annurev.psych.52.1.653>
- Hertwig, R., Benz, B., & Krauss, S. (2008). The conjunction fallacy and the many meanings of and. *Cognition*, *108*, 740–753. <http://dx.doi.org/10.1016/j.cognition.2008.06.008>
- Hertwig, R., & Gigerenzer, G. (1999). The ‘conjunction fallacy’ revisited: How intelligent inferences look like reasoning errors. *Journal of Behavioral Decision Making*, *12*, 275–305. [http://dx.doi.org/10.1002/\(SICI\)1099-0771\(199912\)12:4<275::AID-BDM323>3.0.CO;2-M](http://dx.doi.org/10.1002/(SICI)1099-0771(199912)12:4<275::AID-BDM323>3.0.CO;2-M)
- Hertwig, R., Herzog, S. M., Schooler, L. J., & Reimer, T. (2008). Fluency heuristic: A model of how the mind exploits a by-product of information retrieval. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *34*, 1191–1206. <http://dx.doi.org/10.1037/a0013025>
- Hilbert, M. (2012). Toward a synthesis of cognitive biases: How noisy information processing can bias human decision making. *Psychological Bulletin*, *138*, 211–237. <http://dx.doi.org/10.1037/a0025940>
- Jones, M. N., & Mewhort, D. J. (2007). Representing word meaning and order information in a composite holographic lexicon. *Psychological Review*, *114*, 1–37. <http://dx.doi.org/10.1037/0033-295X.114.1.1>
- Jones, M. N., Willits, J. A., & Dennis, S. (2015). Models of semantic memory. In J. R. Busemeyer & J. T. Townsend (Eds.), *Oxford handbook of mathematical and computational psychology* (pp. 232–254). New York, NY: Oxford University Press.
- Juslin, P., & Persson, M. (2002). PROBABILITIES from EXemplars (PROBEX): A “lazy” algorithm for probabilistic inference from generic knowledge. *Cognitive Science*, *26*, 563–607. http://dx.doi.org/10.1207/s15516709cog2605_2
- Kahneman, D. (2003). A perspective on judgment and choice: Mapping bounded rationality. *American Psychologist*, *58*, 697–720. <http://dx.doi.org/10.1037/0003-066X.58.9.697>
- Kahneman, D., & Frederick, S. (2002). *Representativeness revisited: Attribute substitution in intuitive judgment*. *Heuristics and Biases: The Psychology of Intuitive Judgment*. New York, NY: Cambridge University Press.
- Kahneman, D., & Tversky, A. (1973). On the psychology of prediction. *Psychological Review*, *80*, 237–251. <http://dx.doi.org/10.1037/h0034747>
- Kintsch, W. (2001). Predication. *Cognitive Science*, *25*, 173–202. http://dx.doi.org/10.1207/s15516709cog2502_1
- Kintsch, W. (2014). Similarity as a function of semantic distance and amount of knowledge. *Psychological Review*, *121*, 559–561.
- Koehler, J. J. (1996). The base rate fallacy reconsidered: Descriptive, normative, and methodological challenges. *Behavioral and Brain Sciences*, *19*, 1–17. <http://dx.doi.org/10.1017/S0140525X00041157>
- Kwantes, P. J. (2005). Using context to build semantics. *Psychonomic Bulletin & Review*, *12*, 703–710. <http://dx.doi.org/10.3758/BF03196761>
- Lagnado, D. A., Newell, B. R., Kahan, S., & Shanks, D. R. (2006). Insight and strategy in multiple-cue learning. *Journal of Experimental Psychology: General*, *135*, 162–183. <http://dx.doi.org/10.1037/0096-3445.135.2.162>
- Landauer, T. K., & Dumais, S. T. (1997). A solution to Plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological Review*, *104*, 211–240. <http://dx.doi.org/10.1037/0033-295X.104.2.211>
- Lund, K., & Burgess, C. (1996). Producing high-dimensional semantic spaces from lexical co-occurrence. *Behavior Research Methods, Instruments & Computers*, *28*, 203–208. <http://dx.doi.org/10.3758/BF03204766>
- Manis, M., Dovalina, I., Avis, N. E., & Cardoze, S. (1980). Base rates can affect individual predictions. *Journal of Personality and Social Psychology*, *38*, 231–248. <http://dx.doi.org/10.1037/0022-3514.38.2.231>
- Marewski, J. N., & Schooler, L. J. (2011). Cognitive niches: An ecological model of strategy selection. *Psychological Review*, *118*, 393–437. <http://dx.doi.org/10.1037/a0024143>
- Medin, D. L., & Edelson, S. M. (1988). Problem structure and the use of base-rate information from experience. *Journal of Experimental Psychology: General*, *117*, 68–85. <http://dx.doi.org/10.1037/0096-3445.117.1.68>
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv, 1301.3781*. Retrieved from <https://arxiv.org/abs/1301.3781>
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, & K. Q. Weinberger (Eds.), *Advances in neural information processing systems* (pp. 3111–3119). Redhook, NY: Curran Associates Inc.
- Mitchell, J., & Lapata, M. (2010). Composition in distributional models of semantics. *Cognitive Science*, *34*, 1388–1429. <http://dx.doi.org/10.1111/j.1551-6709.2010.01106.x>
- Morewedge, C. K., & Kahneman, D. (2010). Associative processes in intuitive judgment. *Trends in Cognitive Sciences*, *14*, 435–440. <http://dx.doi.org/10.1016/j.tics.2010.07.004>
- Nilsson, H., Juslin, P., & Winman, A. (2016). Heuristics can produce surprisingly rational probability estimates: Comment on Costello and Watts (2014). *Psychological Review*, *123*, 103–111. <http://dx.doi.org/10.1037/a0039249>
- Nilsson, H., Winman, A., Juslin, P., & Hansson, G. (2009). Linda is not a bearded lady: Configural weighting and adding as the cause of extension errors. *Journal of Experimental Psychology: General*, *138*, 517–534. <http://dx.doi.org/10.1037/a0017351>
- Oaksford, M., & Chater, N. (2007). *Bayesian rationality the probabilistic approach to human reasoning*. New York, NY: Oxford University Press. <http://dx.doi.org/10.1093/acprof:oso/9780198524496.001.0001>
- Pachur, T., & Hertwig, R. (2006). On the psychology of the recognition heuristic: Retrieval primacy as a key determinant of its use. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *32*, 983–1002. <http://dx.doi.org/10.1037/0278-7393.32.5.983>
- Paperno, D., Marelli, M., Tentori, K., & Baroni, M. (2014). Corpus-based estimates of word association predict biases in judgment of word co-occurrence likelihood. *Cognitive Psychology*, *74*, 66–83. <http://dx.doi.org/10.1016/j.cogpsych.2014.07.001>
- Payne, J. W., Bettman, J. R., & Johnson, E. J. (1988). Adaptive strategy selection in decision making. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *14*, 534–552. <http://dx.doi.org/10.1037/0278-7393.14.3.534>
- Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global Vectors for Word Representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*; Vol. 14, pp. 1532–1543). Doha, Qatar.

- Reyna, V. F., Lloyd, F. J., & Brainerd, C. J. (2003). Memory, development, and rationality: An integrative theory of judgment and decision making. In S. Schneider & J. Shanteau (Eds.), *Emerging perspectives on judgment and decision research* (pp. 201–245). New York, NY: Cambridge University Press.
- Rieskamp, J., & Otto, P. E. (2006). SSL: A theory of how people learn to select strategies. *Journal of Experimental Psychology: General*, *135*, 207–236. <http://dx.doi.org/10.1037/0096-3445.135.2.207>
- Schooler, L. J., & Hertwig, R. (2005). How forgetting aids heuristic inference. *Psychological Review*, *112*, 610–628. <http://dx.doi.org/10.1037/0033-295X.112.3.610>
- Shafir, E. B., Smith, E. E., & Osherson, D. N. (1990). Typicality and reasoning fallacies. *Memory & Cognition*, *18*, 229–239. <http://dx.doi.org/10.3758/BF03213877>
- Shah, A. K., & Oppenheimer, D. M. (2008). Heuristics made easy: An effort-reduction framework. *Psychological Bulletin*, *134*, 207–222. <http://dx.doi.org/10.1037/0033-2909.134.2.207>
- Simon, H. A. (1990). Invariants of human behavior. *Annual Review of Psychology*, *41*, 1–20. <http://dx.doi.org/10.1146/annurev.ps.41.020190.000245>
- Slooman, S. A. (1996). The empirical case for two systems of reasoning. *Psychological Bulletin*, *119*, 3–22. <http://dx.doi.org/10.1037/0033-2909.119.1.3>
- Tentori, K., Crupi, V., & Russo, S. (2013). On the determinants of the conjunction fallacy: Probability versus inductive confirmation. *Journal of Experimental Psychology: General*, *142*, 235–255. <http://dx.doi.org/10.1037/a0028770>
- Trueblood, J. S., & Busemeyer, J. R. (2011). A quantum probability account of order effects in inference. *Cognitive Science*, *35*, 1518–1552. <http://dx.doi.org/10.1111/j.1551-6709.2011.01197.x>
- Turney, P. D., & Pantel, P. (2010). From frequency to meaning: Vector space models of semantics. *Journal of Artificial Intelligence Research*, *37*, 141–188.
- Tversky, A. (1977). Features of similarity. *Psychological Review*, *84*, 327–352. <http://dx.doi.org/10.1037/0033-295X.84.4.327>
- Tversky, A., & Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases. *Science*, *185*, 1124–1131. <http://dx.doi.org/10.1126/science.185.4157.1124>
- Tversky, A., & Kahneman, D. (1983). Extensional versus intuitive reasoning: The conjunction fallacy in probability judgment. *Psychological Review*, *90*, 293–315. <http://dx.doi.org/10.1037/0033-295X.90.4.293>
- Tversky, A., & Koehler, D. J. (1994). Support theory: A nonextensional representation of subjective probability. *Psychological Review*, *101*, 547–567. <http://dx.doi.org/10.1037/0033-295X.101.4.547>
- Weber, E. U., & Johnson, E. J. (2009). Mindful judgment and decision making. *Annual Review of psychology*, *60*, 53.
- Yates, J. F., & Carlson, B. W. (1986). Conjunction errors: Evidence for multiple judgment procedures, including “signed summation”. *Organizational Behavior and Human Decision Processes*, *37*, 230–253.

Received February 22, 2016

Revision received August 17, 2016

Accepted October 3, 2016 ■

E-Mail Notification of Your Latest Issue Online!

Would you like to know when the next issue of your favorite APA journal will be available online? This service is now available to you. Sign up at <http://notify.apa.org/> and you will be notified by e-mail when issues of interest to you become available!