

Norms of valence, arousal, and dominance for 13,915 English lemmas

Amy Beth Warriner · Victor Kuperman · Marc Brysbaert

Published online: 13 February 2013
© Psychonomic Society, Inc. 2013

Abstract Information about the affective meanings of words is used by researchers working on emotions and moods, word recognition and memory, and text-based sentiment analysis. Three components of emotions are traditionally distinguished: valence (the pleasantness of a stimulus), arousal (the intensity of emotion provoked by a stimulus), and dominance (the degree of control exerted by a stimulus). Thus far, nearly all research has been based on the ANEW norms collected by Bradley and Lang (1999) for 1,034 words. We extended that database to nearly 14,000 English lemmas, providing researchers with a much richer source of information, including gender, age, and educational differences in emotion norms. As an example of the new possibilities, we included stimuli from nearly all of the category norms (e.g., types of diseases, occupations, and taboo words) collected by Van Overschelde, Rawson, and Dunlosky (Journal of Memory and Language 50:289–335, 2004), making it possible to include affect in studies of semantic memory.

Keywords Emotion · Semantics · Gender differences · Age differences · Crowdsourcing

Emotional ratings of words are in high demand because they are used in at least four lines of research. The first of these lines concerns research on the emotions themselves: the ways

in which they are produced and perceived, their internal structure, and the consequences that they have for human behavior. For instance, Verona, Sprague, and Sadeh (2012) used emotionally neutral and negative words in an experiment comparing the responses of offenders without a personality disorder to those of offenders with an antisocial personality disorder who either did or did not have additional psychopathic traits.

The second line of research deals with the impact that emotional features have on the processing and memory of words. Kousta, Vinson, and Vigliocco (2009) found that participants responded faster to positive and negative words than to neutral words in a lexical-decision experiment, a finding later replicated by Scott, O'Donnell, and Sereno (2012) in sentence reading. According to Kousta, Vigliocco, Vinson, Andrews, and Del Campo (2011), emotion is particularly important in the semantic representations of abstract words. In other research, Fraga, Piñeiro, Acuña-Fariña, Redondo, and García-Orza (2012) reported that emotional words are more likely to be used as attachment sites for relative clauses in sentences such as “Someone shot the servant of the actress who. . . .”

A third approach uses emotional ratings of words to estimate the sentiments expressed by entire messages or texts. Leveau, Jhean-Larose, Denhière, and Nguyen (2012), for instance, wrote a computer program to estimate the valence and arousal evoked by texts on the basis of word measures (see also Liu, 2012).

Finally, emotional ratings of words are used to automatically estimate the emotional values of new words by comparing them to those of validated words. Bestgen and Vincze (2012) gauged the affective values of 17,350 words by using the rated values of words that were semantically related.

So far, nearly all studies have been based on Bradley and Lang's (1999) Affective Norms for English Words (ANEW) or on translated versions (for exceptions, see Kloumann, Danforth, Harris, Bliss, & Dodds, 2012; Mohammad &

Electronic supplementary material The online version of this article (doi:10.3758/s13428-012-0314-x) contains supplementary material, which is available to authorized users.

A. B. Warriner · V. Kuperman (✉)
Department of Linguistics and Languages, McMaster University,
Togo Salmon Hall 626, 1280 Main Street West,
Hamilton, Ontario L8S 4M2, Canada
e-mail: vickup@mcmaster.ca

M. Brysbaert
Ghent University, Ghent, Belgium

Turney, 2010). These norms include ratings for 1,034 words. Three types of ratings were carried out, in line with Osgood, Suci, and Tannenbaum's (1957) theory of emotions. The first, and most important, type of ratings concerns the valence (or pleasantness) of the emotions invoked by a word, going from *unhappy* to *happy*. The second addresses the degree of arousal evoked by a word, and the third dimension refers to the dominance/power of the word—the extent to which the word denotes something that is weak/submissive or strong/dominant.

The number of words covered by the ANEW norms appeared sufficient for use in small-scale factorial experiments. In these experiments, a limited number of stimuli would be selected that varied on one dimension (e.g., valence) and were matched on other variables (e.g., arousal, word frequency, and word length). However, the number of words in this set is prohibitively small for the large-scale megastudies that are currently emerging in psycholinguistics. In these studies (e.g., Balota et al., 2007; Ferrand et al., 2010; Keuleers, Brysbaert, & New, 2010; Keuleers, Lacey, Rastle, & Brysbaert, 2012), regression analyses of thousands of words are used to disentangle the influences on word recognition. The ANEW norms are also limited as input for computer algorithms that gauge the sentiment of a message/text or the emotional values of nonrated words.

Given the ease with which word norms can be collected nowadays, we decided to collect affective ratings for a majority of the well-known English content words (a total of 13,915). Because it would be expected that the emotional values would generalize to inflected forms (e.g., *sings*, *sang*, *sung*, and *singing* for the verb lemma *sing*), **we only included lemmas (the base forms of words—i.e., the ones used as entries in dictionaries)**. Our sample of words (see below for the selection criteria) substantially covers the word stock of the English language and forms a solid foundation from which to automatically derive the values of the remaining words (Bestgen & Vincze, 2012).

Method

Stimuli

The words included in our stimulus set were compiled from three sources: Bradley and Lang's (1999) ANEW database, Van Overschelde, Rawson, and Dunlosky's (2004) category norms, and the SUBTLEX-US corpus (Brysbaert & New, 2009). Our final set included 1,029 of the 1,034 words from ANEW (five were lost due to programmatic error) and 1,060 of the participant-generated responses to 60 of the 70 category names included in the category norm study **(we did not include a few categories, such as units of time and distance or types of fish)**. The remaining words were selected from the

list of 30,000 lemmas for which Kuperman, Stadthagen-Gonzalez, and Brysbaert (2012) collected age-of-acquisition ratings. This list contains the content lemmas (nouns, verbs, and adjectives) from the 50-million-token SUBTLEX-US subtitle corpus. We only selected the highest-frequency words known by 70 % or more of the participants in Kuperman et al., given that affective ratings are less valid/useful for words that are not known to most participants. Our final set included 13,915 words, of which 22.5 % are most often used as adjectives (Brysbaert, New, & Keuleers, 2012), 63.5 % as nouns, 12.6 % as verbs, and 1.4 % as other or unspecified parts of speech. The mean word frequency of the set was 1,056 ($SD = 8,464$, range = 1 to 314,232, median = 87) in the 50-million-token SUBTLEX-US corpus; 152 words, or 1 %, had no frequency data. For each word in our set, we collected ratings on three dimensions using a 9-point scale.

The stimuli were distributed over 43 lists containing 346 to 350 words each. Each list consisted of 10 calibrator words, 40 control words from ANEW, and a randomized selection of non-ANEW words. The calibrator words were drawn from ANEW and were chosen separately for each of the three dimensions, with the goal of giving participants a sense of the entire range of the stimuli that they would encounter.¹ Participants always saw these calibrator words first. The remaining ANEW words were divided into sets of 40 and served as controls for the estimation of correlations between our data and the ANEW norms. This meant that a selection of these words appeared in more than one list and that the lists used for each of the three dimensions were mostly, but not completely, identical. The control words and the non-ANEW words were randomly mixed together in each list. Once lists were created, the words in each one were always presented in a fixed order following the calibrator words.

Data collection

Participants were recruited via the Amazon Mechanical Turk crowdsourcing website. Responders were restricted to those who self-identified as being current residents of the US and who completed any given list only once. This completion of a single list by a given participant will henceforth be referred to as an *assignment*. Each assignment involved rating words on a

¹ The calibrator words for the respective dimensions were as follows (in increasing order of ratings): *Valence*: "jail" (1.91), "invader" (2.23), "insecure" (2.30), "industry" (5.07), "icebox" (5.67), "hat" (5.69), "grin" (7.66), "kitten" (7.58), "joke" (7.88), and "free" (8.25). *Arousal*: "statue" (2.82), "rock" (3.14), "sad" (3.49), "cat" (4.50), "curious" (5.74), "robber" (6.20), "shotgun" (6.55), "assault" (6.80), "thrill" (7.19), and "sex" (7.60). *Dominance*: "lightning" (4.00), "mildew" (4.19), "waterfall" (5.34), "wealthy" (6.11), "lighthouse" (6.24), "honey" (6.39), "treat" (6.66), "mighty" (6.85), "admired" (6.94), and "liberty" (7.04).

single dimension only, in contrast to the ANEW study, for which participants rated each word on all three dimensions. The instructions given were minor variations on the instructions in the ANEW project, and are given below, with the respective changes to the wording for the separate dimensions indicated in square brackets.

You are invited to take part in the study that is investigating emotion, and concerns how people respond to different types of words. You will use a scale to rate how you felt while reading each word. There will be approximately 350 words. The scale ranges from 1 (*happy [excited; controlled]*) to 9 (*unhappy [calm; in control]*). At one extreme of this scale, you are happy, pleased, satisfied, contented, hopeful [stimulated, excited, frenzied, jittery, wide-awake, or aroused; controlled, influenced, cared-for, awed, submissive, or guided]. When you feel completely happy [aroused; controlled] you should indicate this by choosing rating 1. The other end of the scale is when you feel completely unhappy, annoyed, unsatisfied, melancholic, despaired, or bored [relaxed, calm, sluggish, dull, sleepy, or unaroused; in control, influential, important, dominant, autonomous, or controlling]. You can indicate feeling completely unhappy [calm; in control] by selecting 9. The numbers also allow you to describe intermediate feelings of pleasure [calmness/arousal; in/under control], by selecting any of the other feelings. If you feel completely neutral, neither happy nor sad [not excited nor at all calm; neither in control nor controlled], select the middle of the scale (rating 5).

Please work at a rapid pace and don't spend too much time thinking about each word. Rather, make your ratings based on your first and immediate reaction as you read each word.

On average, assignments were completed in approximately 14 min. Participants received 75 cents per completed assignment. After reading an informational consent statement and the instructions, participants were asked to indicate their age, gender, first language(s), country/state resided in most between birth and age 7, and educational level. Subsequently, they were reminded of the scale anchors and presented with a scrollable page in which all words in the list were shown to the left of nine numbered radio buttons. Although we did not incorporate the Self-Assessment Manikins (SAM) that were used in the ANEW study, we did anchor our scales in the same direction, with valence ranging from *happy* to *unhappy*, arousal from *excited* to *calm*, and dominance from *controlled* to *in control*. In the [Results and Discussion](#) section, we show that our numerical ratings correlated highly with the SAM ratings from ANEW, demonstrating that the methods are roughly equivalent. Once finished, participants clicked "Submit" to complete the study.

Lists were initially presented to 20 respondents each. However, missing values due to subsequent exclusion criteria resulted in some words having fewer than 18 valid ratings. Several of the lists were reposted until the vast majority of the words had reached at least this threshold. Data collection began on March 14, 2012, and was completed May 30, 2012.

Results and discussion

Data trimming

Altogether, 1,085,998 ratings were collected across all three dimensions. Around 3 % of the data were removed due to missing responses, lack of variability in responses (i.e., providing the same rating for all words in the list), or the completion of fewer than 100 ratings per assignment. The valence and arousal ratings were reversed post-hoc to maintain a more intuitive low-to-high scale (e.g., *sad* to *happy* rather than *happy* to *sad*) across all three dimensions. Means and standard deviations were calculated for each word. Ratings in assignments with negative correlations between a given participant's rating and the mean for that word were reversed (9 %). This was done on the basis of both empirical evidence that higher numbers intuitively go with positive anchors (Rammstedt & Krebs, 2007) and an examination of these participants' responses, which revealed unintuitive answers (e.g., indicating that negative words such as "jail" made them very happy). Any remaining assignments with ratings that correlated with the mean ratings per items at less than .10 were removed, and the means and standard deviations were recalculated. The final data set consisted of 303,539 observations for valence (95 % of the original data pool), 339,323 observations for arousal (89 % of the original data pool), and 281,735 observations for dominance (74 % of the original data pool). A total of 1,827 responders contributed to this final data set, with 362 of them completing assignments for two or more dimensions. A total of 144 participants completed two or more assignments within a single dimension.

For valence, 51 words received fewer than 18 (but more than 15) valid ratings. For arousal, 128 words had a total number of ratings in that range. For dominance, 564 words had a total of either 16 or 17 ratings, and 17 words had 14 or 15 ratings each. For all three dimensions, more than 87 % of the words had between 18 and 30 ratings per word. A total of 50 words in each dimension received more than 70 ratings each, due to the doubling up of ANEW words and the rerunning of lists. To illustrate how our data enriches the set of words available in ANEW, Table 1 provides examples of words that are not included in the ANEW list and that show very high or very low ratings in one of the three dimensions.

Demographics

Of the 1,827 valid responders, approximately 60 % were female in all three cases (419 valence, 448 arousal, and 505 dominance). Their ages ranged from 16 to 87 years, with 11 % being 20 years old or younger; 45 % from 21 to 30; 21 % from 31 to 40; 11 % from 41 to 49; and 12 % age 50 or

Table 1 Words at the extremes of each dimension that were not included in ANEW

	Valence		Arousal		Dominance	
Lowest	pedophile	1.26	grain	1.60	dementia	1.68
	rapist	1.30	dull	1.67	Alzheimer's	2.00
	AIDS	1.33	calm	1.67	lobotomy	2.00
	leukemia	1.47	librarian	1.75	earthquake	2.14
	molester	1.48	soothing	1.91	uncontrollable	2.18
	murder	1.48	scene	1.95	rapist	2.21
Highest	excited	8.11	motherfucker	7.33	rejoice	7.68
	sunshine	8.14	erection	7.37	successful	7.71
	relaxing	8.19	terrorism	7.42	smile	7.72
	lovable	8.26	lover	7.45	completion	7.73
	fantastic	8.36	rampage	7.57	self	7.74
	happiness	8.48	insanity	7.79	incredible	7.74

older. Of the participants, 24 (3.3 %), 32 (4.3 %), and 23 (2.7 %) for the valence, arousal, and dominance dimensions, respectively, reported a native language other than English, while 10 (1.4 %), 12 (1.6 %), and 12 (1.4 %) participants, respectively, reported more than one native language, including English. Table 2 shows the numbers of participants at each of the seven possible education levels. Most had some college or a bachelor's degree.

Descriptive statistics

Table 3 reports descriptive statistics for the three distributions of ratings. The distributions of both valence and dominance ratings are negatively skewed ($G_1 = -.28$ and $-.23$, respectively), with 55 % of the words rated above the median of the rating scale for both dimensions (see Fig. 1). The Mann–Whitney one-sample median test indicated that the medians of both the valence and dominance

distributions were not significantly different from rating 5, which is the median of the scales (both $ps > .1$). **The tendency for more words to make people feel happy and in control goes along with numerous former findings of positivity biases in English and other languages (see Augustine, Mehl, & Larsen, 2011, and Kloumann et al., 2012).** The *positivity bias*—or the prevalence of positive word types in English books, Twitter messages, music lyrics, and other genres of texts—is argued to reflect the preference of humankind for pro-social and benevolent communication. Arousal, on the other hand, is positively skewed ($G_1 = .47$), meaning that only a relatively small proportion of words (20 % above a rating of 5) made people feel excited.

Ratings of valence were relatively consistent across participants, while arousal and dominance were much more variable. This is indicated by the difference between the average standard deviations of the dimensions: 1.68 for valence, but 2.30 and 2.16 for arousal and dominance, respectively. In addition, the split-half reliabilities were .914 for valence, .689 for arousal, and .770 for dominance; see below for other examples of a higher variability of dominance and arousal ratings. Figure 2a–c show, for the three emotional dimensions, the means of the ratings for each word plotted against their standard deviations, with each scatterplot's smoother lowess line demonstrating the

Table 2 Reported education levels within each dimension

Education Level	Number of Participants		
	Valence (%)	Arousal (%)	Dominance (%)
Some high school	28 (4)	32 (4)	28 (3)
High school graduate	96 (13)	98 (13)	117 (14)
Some college–No degree	237 (33)	252 (34)	298 (35)
Associates degree	82 (11)	79 (11)	93 (11)
Bachelors degree	212 (29)	222 (30)	218 (26)
Masters degree	55 (8)	53 (7)	78 (9)
Doctorate	13 (2)	9 (1)	13 (2)
Total	723	745	845

The numbers across all three columns add up to more than 1,827, as some people contributed to more than one dimension

Table 3 Descriptive statistics for the distribution of each dimensions, including the number of participants (N), number of observations, average mean, and average SD

	N	# of Obs	Mean	Avg SD
Valence	723	303,539	5.06	1.68
Arousal	745	339,323	4.21	2.30
Dominance	845	281,735	5.18	2.16

Table 5 Correlations of present ratings with similar studies across languages

Data Set				Correlations		
Source	Language	<i>N</i> (source)	<i>N</i> (overlap)	Valence	Arousal	Dominance
a	English	1,040	1,029	.953	.759	.795
b	Dutch	4,299	3,701	.847	.575	N/A
c	Spanish	1,034	1,023	.924	.692	.833
d	Portuguese	1,040	1,023	.924	.635	.774
e	Finnish	213	203	.956	N/A	N/A
f	English	10,222	4,504	.919	N/A	N/A

Sources: a, Bradley & Lang (1999); b, Moors et al. (in press)—English glosses; c, Redondo, Fraga, Padrón, & Comesaña, (2007)—English glosses; d, Soares, Comesaña, Pinheiro, Simões, & Frade (2012)—English glosses; e, Eilola & Havelka (2010)—English glosses; f, Kloumann, Danforth, Harris, Bliss, & Dodds (2012). All studies except Moors et al. (in press) utilized a nine-point scale in acquiring their ratings. Moors et al. (in press) used a seven-point scale

Table 6 Correlations between emotional dimensions and semantic variables reported in prior studies [degrees of freedom are based on the numbers of data points reported as *N* (Overlap)]

Source	Measure	<i>N</i> (Source)	<i>N</i> (Overlap)	Valence	Arousal	Dominance
a	Imageability	5,988	5,125	.161	−.012	.031
b	Imageability	326	318	−.037	.099	−.160
	Concreteness	326	318	.109	−.244	−.019
	Context Avail.	326	318	.196	−.147	.044
c	Concreteness	1,944	1,567	.105	−.258	.009
d	Imageability	3,394	2,906	.152	−.045	.006
	Familiarity	3,394	2,906	.206	−.028	.215
e	AoA ¹	30,121	13,709	−.233	−.062	−.187
	% Known ²	30,121	13,709	.094	.078	.103
f	Sensory Exp.	5,857	5,007	.067	.228	−.044
g	Body–Object	1,618	1,398	.203	−.143	.172
h	Familiarity	559	503	.272	−.193	.329
	Pain	559	503	−.456	.579	−.343
	Smell	559	503	.139	.052	−.043
	Color	559	503	.401	.052	.081
	Taste	559	503	.309	−.102	.084
	Sound	559	503	−.176	.407	−.286
	Grasp	559	503	.024	−.121	.252
	Motion	559	503	−.113	.328	−.328
i	Sound	1,402	1,283	−.04	.311	−.121
	Color	1,402	1,283	.322	−.072	.100
	Manipulation	1,402	1,283	.070	.026	.255
	Motion	1,402	1,283	.011	.335	−.140
	Emotion	1,402	1,283	.902	−.206	.658
j	Log Frequency ³	74,286	13,763	.182	−.033	.167

¹ AoA, age of acquisition. ² The overlapping words in this study represent a biased sample, due to the fact that words in the present study were restricted to only include words that were known by 70 % or more participants in the studies cited here. ³ Since we chose words to fill our quota that were higher in frequency, the overlap here is also biased toward the upper range

Sources: a, Cortese & Fugett (2004) and Schock, Cortese, & Khanna (2012); b, Altarriba, Bauer, & Benvenuto (1999); c, Gilhooly & Logie (1980); d, Stadthagen-Gonzalez & Davis (2006); e, Kuperman, Stadthagen-Gonzalez, & Brysbaert (2012); f, extended data set of Juhasz & Yap (in press) and Juhasz, Yap, Dicke, Taylor, & Gullick (2011); g, Tillotson, Siakaluk, & Pexman (2008); h, Amsel, Urbach, & Kutas (2012); i, Medler, Arnoldussen, Binder, & Seidenberg (2005); j, Brysbaert & New (2009)

including the ANEW set from which we drew our control words. The correlations are listed in Table 5.

Valence appears to generalize very well across studies and languages, as evidenced by high correlations. Both arousal and dominance show more variability across languages and studies, as reflected in the lower correlations. Note that these studies themselves (those that have reported the information—i.e., c, d, and e) also found a lower correlation between their arousal and dominance ratings and the arousal and dominance ratings reported in other studies (arousal range = .65 to .75; dominance range = .72 to .73). Importantly, however, cross-linguistic correlations were stronger (the range of Pearson's r for arousal was .575–.759) than those between gender, age, and education groups within our study (the range of Pearson's r was .467–.516), see Table 8 below. This observation clearly indicates the validity of using emotional ratings to English glosses of words in a language that does not have an extensive set of ratings at the researcher's disposal. This seems to be more the case for valence and dominance than for arousal.

Correlations with lexical properties

As is known for other subjective ratings of lexical properties (cf. Baayen, Feldman, & Schreuder, 2006), judgments of the

emotional impact of a word are likely to be affected by other aspects of the word's meaning. Table 6 reports correlations of valence, arousal, and dominance with a range of available semantic variables. In the remainder of the article, words, rather than the trial-level data, were chosen as units of the correlational analyses.

Most of the correlations that the emotional ratings show with other semantic properties are weak to moderate (Cohen, 1992), with the exception of correlations with variables that directly tap into emotional states (h and i in Table 6). Specifically, words that make people happy are easier to picture [$r(5123) = .161, p < .001$] and more concrete [$r(1565) = .105, p < .001$], familiar [$r(2904) = .206, p < .001$], context rich [$r(316) = .196, p < .001$], and easy to interact with [$r(1396) = .203, p < .001$], are of high frequency [$r(13763) = .182, p < .001$], and are learned at an early age [$r(13707) = -.233, p < .001$]. They are also associated with low pain [$r(501) = -.456, p < .001$], intense smell [$r(501) = .139, p < .01$], vivid color [$r(1281) = .322, p < .001$], pleasant taste [$r(501) = .309, p < .001$], quiet sounds [$r(501) = -.176, p < .001$], and stillness [$r(501) = -.113, p < .05$]. Virtually all of these properties are also associated with words that make people feel in control; that is, they correlate in the same way with dominance ratings.

Words that make people feel excited are more ambiguous [$r(1565) = -.258, p < .001$], unfamiliar [$r(501) = -.193, p < .001$],

Fig. 4 Relationships between the three dimensions and age of acquisition, word frequency, imageability, and sensory experience ratings, presented as scatterplot smoother lowess trend lines

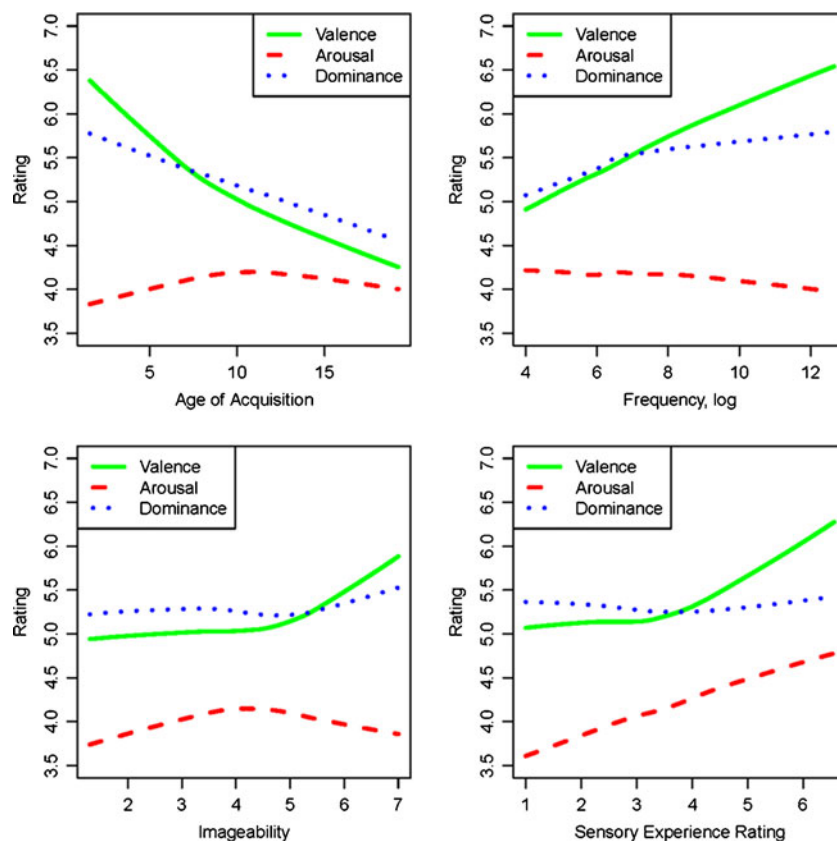


Table 7 Group differences in emotional dimensions

Measure	<i>N</i>	# of Obs	Mean	Avg <i>SD</i>	<i>N</i>	# of Obs	Mean	Avg <i>SD</i>	<i>p</i>
Male					Female				
Valence	301	116,819 (38 %)	5.13	1.60	419	184,636 (61 %)	5.00	1.64	<.001
Arousal	291	119,658 (37 %)	4.38	2.27	448	197,648 (62 %)	4.10	2.28	<.001
Dominance	336	149,329 (44 %)	4.83	2.15	505	188,433 (55 %)	4.81	2.13	n.s.
Old					Young				
Valence	346	158,067 (52 %)	5.04	1.61	382	147,892 (48 %)	5.10	1.68	<.001
Arousal	373	174,402 (54 %)	4.13	2.27	374	146,021 (46 %)	4.31	2.31	<.001
Dominance	384	153,581 (45 %)	4.80	2.04	464	187,137 (55 %)	4.88	2.17	<.001
High Education					Low Education				
Valence	362	136,280 (45 %)	5.10	1.57	361	167,259 (55 %)	5.04	1.70	<.05
Arousal	363	142,151 (45 %)	4.28	2.17	382	177,213 (55 %)	4.14	2.33	<.001
Dominance	402	154,590 (46 %)	5.17	2.02	443	184,733 (54 %)	5.20	2.22	<.05

Reported are the numbers of raters (*N*), numbers of observations (# of Obs), and percentages of total observations in each group (in parentheses), the group means and the average standard deviations, and, in the last column, the *p* value of a two-tailed independent *t* test comparing the group means. The numbers of observations do not always equal 100 % because a small number of participants declined to answer the relevant demographic questions.

context impoverished [$r(316) = -.147, p < .01$], and difficult to interact with [$r(1396) = -.143, p < .001$]. They are also associated with strong general sensory experience [$r(5005) = .228, p < .001$], specifically with high pain [$r(501) = .579, p < .001$], unpleasant taste [$r(501) = -.102, p < .05$], intense sounds [$r(501) = .407, p < .001$], motion [$r(1281) = .335, p < .001$], and an inability to be grasped [$r(501) = -.121, p < .01$].

As correlations do not reveal the form of the functional relationships, Fig. 4 below zooms in on functional relationships between the three emotional dimensions and selected semantic properties of interest.

The top left panel of Fig. 4 reveals that early words are maximally positive, strong, and calm. Words become more negative and weak (controlled by) on average as the age of acquisition increases. The peak of arousal is reached in the words learned around the age of 10, while later-acquired words are less exciting. It is tempting to interpret these results as an average developmental timeline of vocabulary acquisition in North American children, with (a) earliest happy and calm words learned in a risk-averse environment protecting a child from negativity and excitement, and (b) excitable words like sexual terms, taboo words, and swear words learned in early school age. Yet it is more likely that the age-of-acquisition patterns of emotional words are at least partly due to how often they occur in English, and thus how likely children are to encounter and learn them early. The top right of Fig. 4 demonstrates that the more frequent a word is, the happier, stronger, and calmer it tends to be. The observed linear relationship between log frequency of occurrence and valence is reasonably strong: The Pearson's correlation coefficient is .18, and the increase in valence between the least and most frequent words is on the order of

two points on the 9-point scale. This corroborates the finding of Garcia, Garas, and Schweitzer (2012) and runs counter to the claim of Kloumann et al. (2012) that the positivity bias in English words is only observed in word types (there are more positive than negative words) and that the correlations between frequency and valence, if any, are corpus-specific and small. The discrepancy may be due to the much broader range of frequency that we consider here, with 14,000 words from the top of the frequency list rather than 5,000 words in each of the corpora considered by Kloumann et al. We leave the verification of the positivity bias over a broader frequency range to further research.

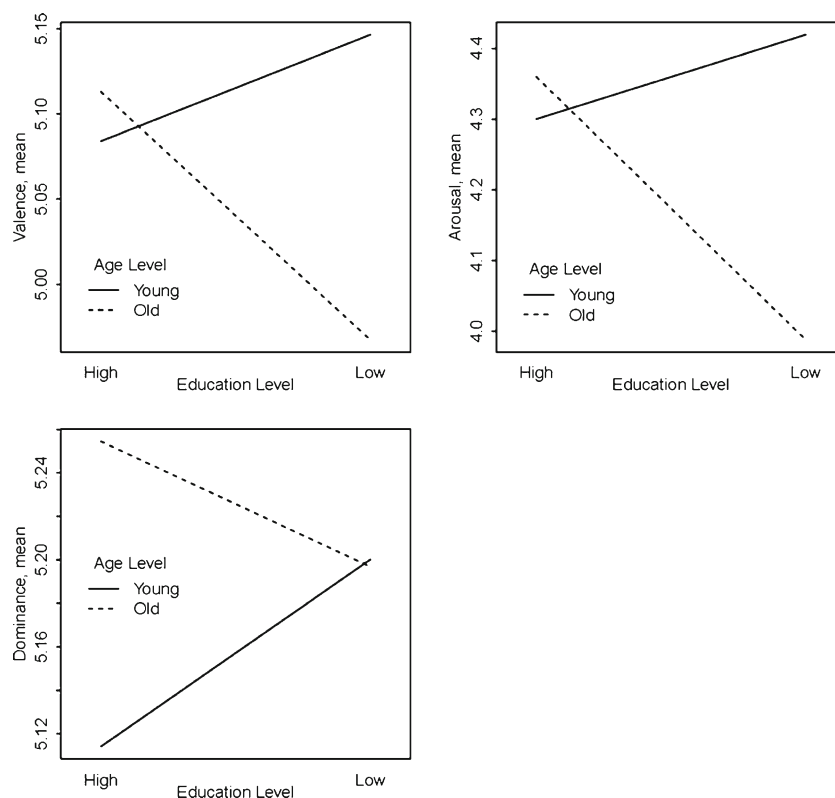
Only highly imageable words are emotionally colored (Fig. 4, bottom left): As imageability increases from rating 5 on the 7-point scale, words become more positive and strong (in control). Again, arousal is distinct from this pattern: Words that are hardly imageable at all or very imageable are calm, while those in the middle of the imageability range increase excitement.

The increasing strength of a sensory experience (Fig. 4, bottom right) varies strongly with arousal: The more tangible the word is, the more exciting it is. This suggests that abstract notions are less powerful in agitating human readers than are material objects. The functional relationship with

Table 8 Correlations between groups

	Valence	Arousal	Dominance
Male and female	.789	.516	.593
Old and young	.818	.500	.591
High and low education	.831	.467	.608

Fig. 5 Interactions between dichotomized education and age levels for all three dimensions. All interactions are significant at $p < .001$



valence is only observed in the top half of the sensory experience range: More tangible words induce increasingly positive emotions. No reliable relationship is observed between sensory experience ratings and dominance.

Interactions between demographics and ratings

Participants were naturally divided into two genders. In addition, we divided them into two age ranges using the median split—younger (less than 30) and older (30 or greater). We also dichotomized education level into higher (those who had an associate's degree or greater) and lower (some college or less). All three dimensions showed slightly but significantly higher average ratings for younger versus older and for lower education versus higher education. Also, males gave slightly but reliably higher ratings in all dimensions than did females. Separate independent t tests showed that this difference was significant for valence and arousal, but not for dominance. The means, standard deviations, and independent t test significance levels of each group division are listed in Table 7.

Table 8 reports correlations between groups of participants and demonstrates substantial variability in the ratings that they provided: As with the overall data in Table 5, arousal and dominance elicited less agreement in judgments than did valence.

We ran a series of multiple regressions looking at age, gender, and education (all dichotomized as described above)

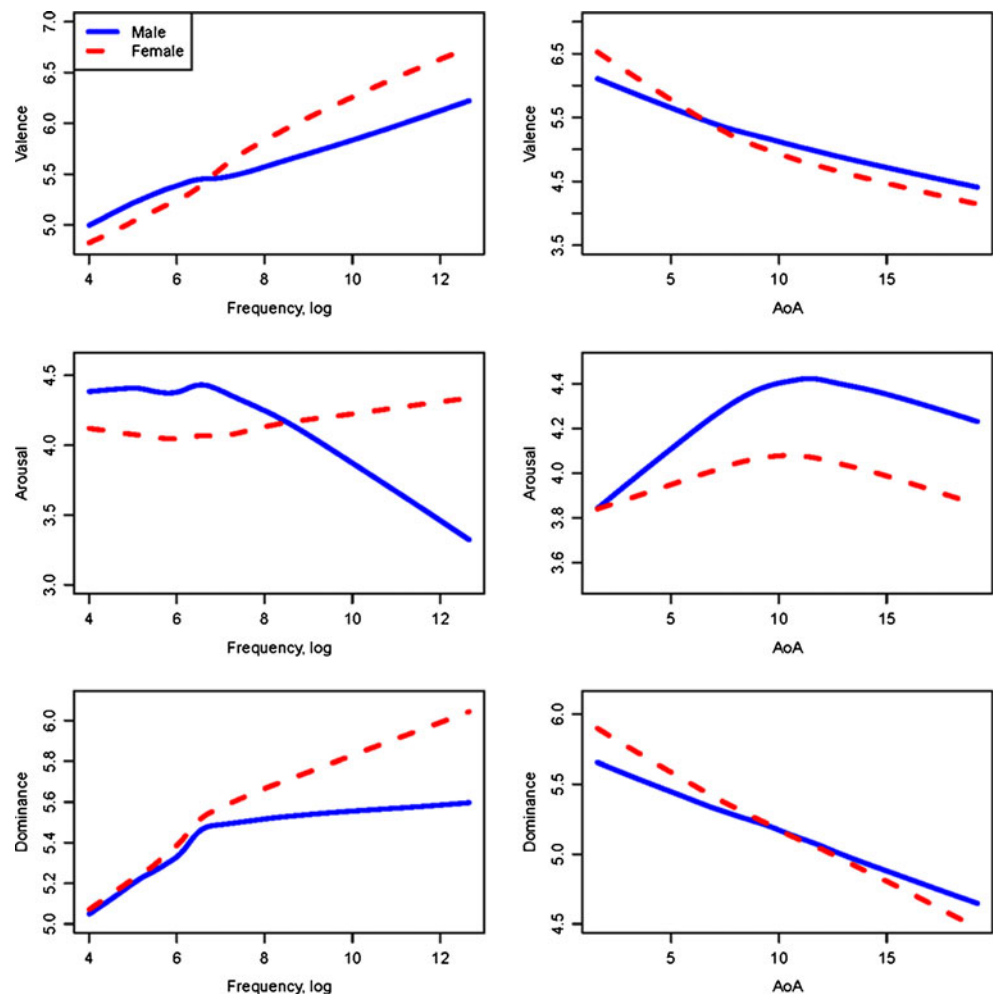
as predictors. All main effects were significant at $p < .001$, and each variable made a unique contribution to the variance in the collected ratings. In addition, most of the two- and three-way interactions for all three dimensions were significant, likely due to the large number of data points available. However, the actual ranges of the effects tended to be small. One exception was the interaction between age and education level for all three dimension (see Fig. 5). For valence and arousal, highly educated people rated words similarly, regardless of age. For those with less education, age strongly affected ratings, with the younger group providing higher ratings, on average, than did the older. For dominance, the opposite pattern held: Age affected those in the higher education group, with older participants providing higher ratings than younger ones, but age did not have an effect in the lower education group.

Gender differences

In what follows, we concentrate on gender differences. Effects of well-established lexical properties on emotion norms varied by gender. Figure 6 presents interactions of gender with frequency of occurrence and age of acquisition as predictors of emotional ratings. All interactions reached significance in multiple regression models, with each set of ratings treated separately as a dependent variable, all $ps < .01$.

The interactions revealed that female raters provided more extreme negative/weak ratings for the lowest-

Fig. 6 Interactions of gender with frequency (left) and age of acquisition (AoA, right) as predictors of mean ratings of valence (top), arousal (middle), and dominance (bottom). Interactions are presented with gender-specific lowest trend lines



frequency words, and more extreme positive/strong ratings for higher-frequency words, yielding a broader range of values for both valence and dominance. The same holds for the more extreme ratings given by females to earliest- and latest-learned words, as compared to males.

Quite the opposite pattern was observed in the ratings of arousal (Fig. 6, middle row). Female raters showed a weak relationship between either frequency or age of acquisition and arousal, with slightly higher arousal words in the higher-frequency band and in the mid-range of age of acquisition. Conversely, male raters revealed a strong tendency to find higher-frequency and earlier-learned words as being less exciting than relatively late and infrequent words.

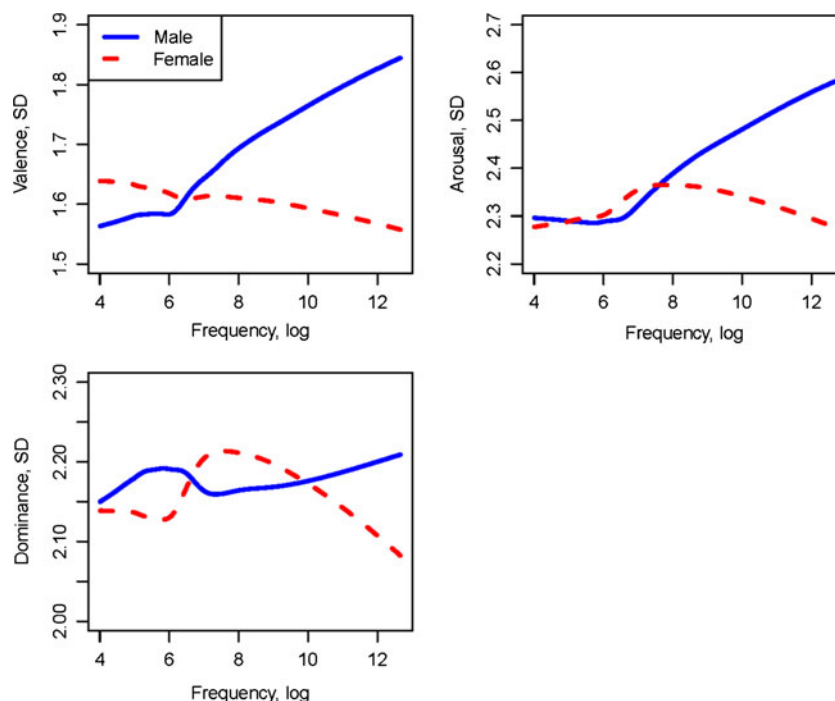
Variability in ratings also varied by gender, see Fig. 7. Male raters disagreed increasingly more on all ratings to higher-frequency words, while variance in ratings by female participants was increasingly attenuated with an increase in word frequency.

While pinning down the origin of these differences will be an issue for further investigation, here we note the necessity for research into emotion words to take into account these interactions as potential sources of systematic error.

Semantic categories

An interesting aspect of emotional ratings is their use to quantify attitudes and opinions toward physical, psychological, and social phenomena either in the population at large or in specific target groups. We showcase here emotional ratings to the semantic categories of disease (Fig. 8) and occupation (Fig. 9), based on Van Overschelde et al.'s (2004) category norms, with occasional additions of semantically similar words. As Fig. 8 suggests, all diseases are rated as words evoking negative feelings, high arousal, and feelings of being controlled; that is, all ratings were below the median of valence/dominance and above the median of arousal in the entire data set (shown as a dotted line). Sexually transmitted diseases were judged as being among the most negative and the most anxiety-provoking entries in the subset. This is generally in line with surveys of attitudes that list sexually transmitted diseases as being among the most stigmatized medical conditions (e.g., Brems, Johnson, Warner, & Roberts, 2010). The most feared medical conditions—cancer, Alzheimer's, heart disease, and stroke (listed by decreasing percentages of respondents who feared

Fig. 7 Interactions of gender with frequency as a predictor of the standard deviations of ratings of valence (top left), arousal (top right), and dominance (bottom left). Interactions are presented with gender-specific lowess trend lines



them; MetLife Foundation, 2011; YouGov, 2011)—are also among the most negative, the least controllable, and the most anxiety-provoking diseases.

Ratings of valence to occupations revealed that the best-paying professions in the list were judged as being the most

negative, below the median in the overall data set: compare “lawyer,” “dentist,” and “manager.” The correlation between average income, as reported by the Bureau of Labor Statistics (2011), and mean valence is indeed negative, but it does not reach significance ($r = -.167$, $p = .434$), possibly

Fig. 8 Ratings of words denoting disease. Dotted lines represent the median ratings of the respective emotional dimensions across the entire data set

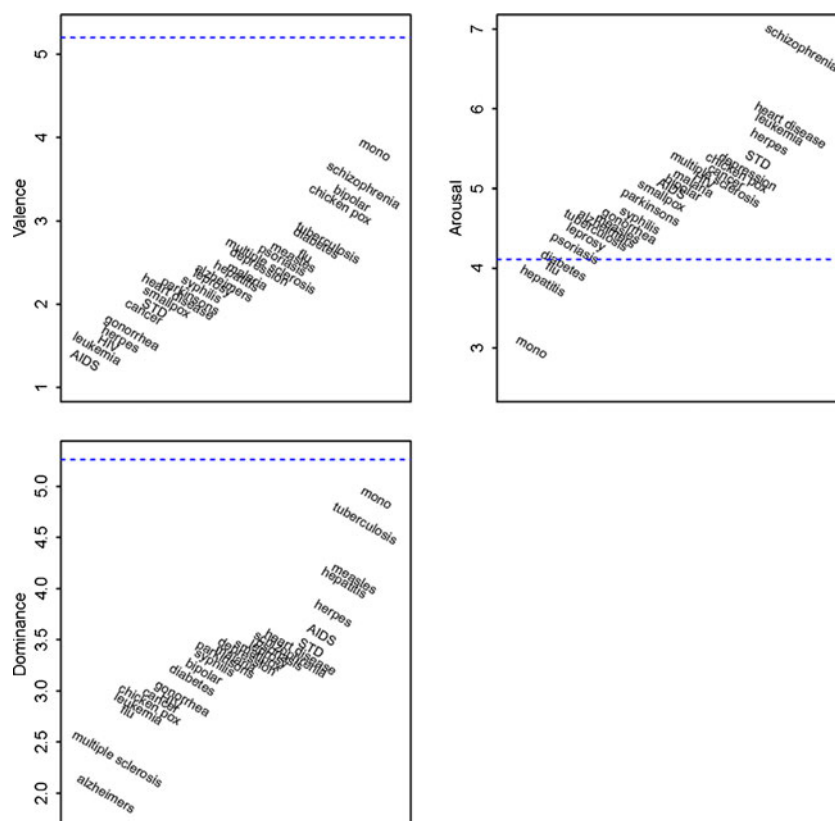


Fig. 9 Ratings of words denoting occupations. Dotted lines represent the median ratings of the respective emotional dimensions across the entire data set

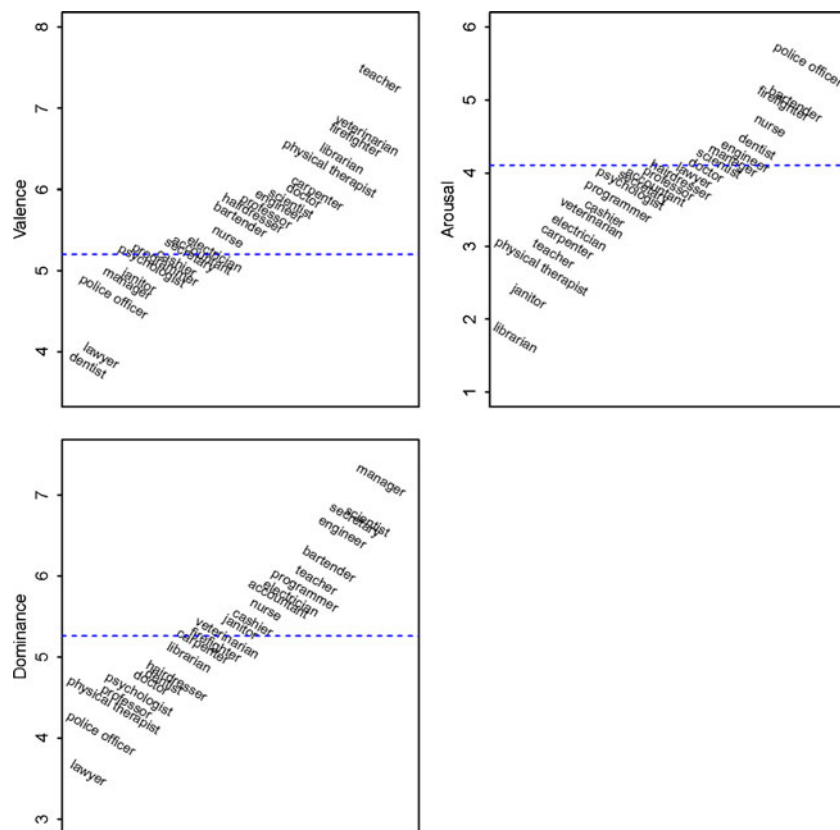


Fig. 10 Gender differences in ratings for weapon-related words

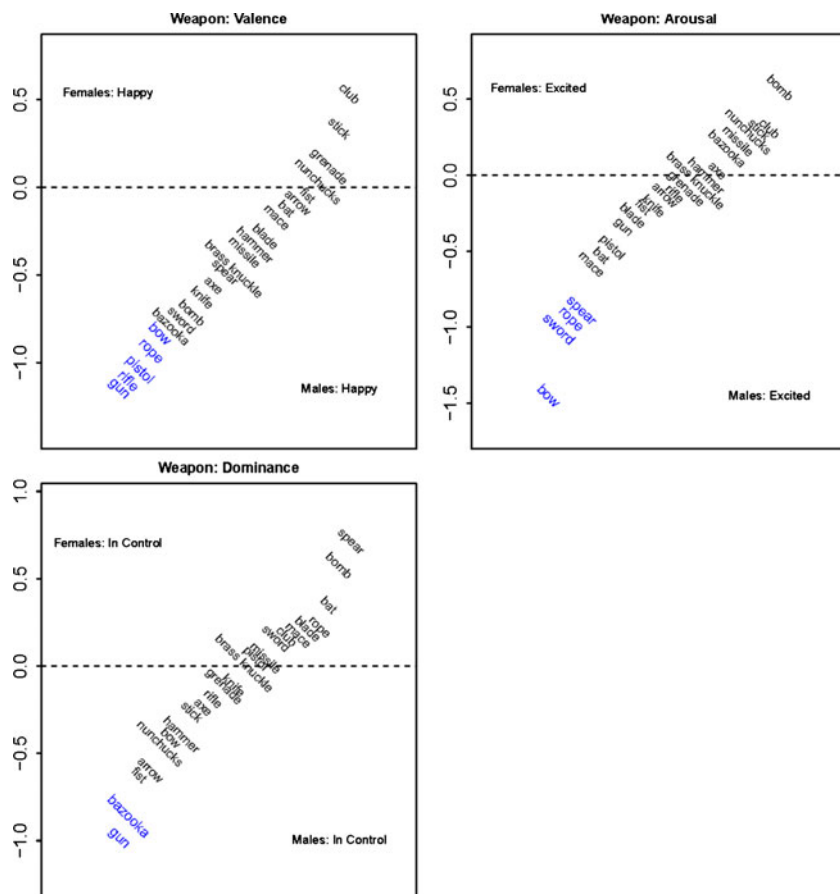
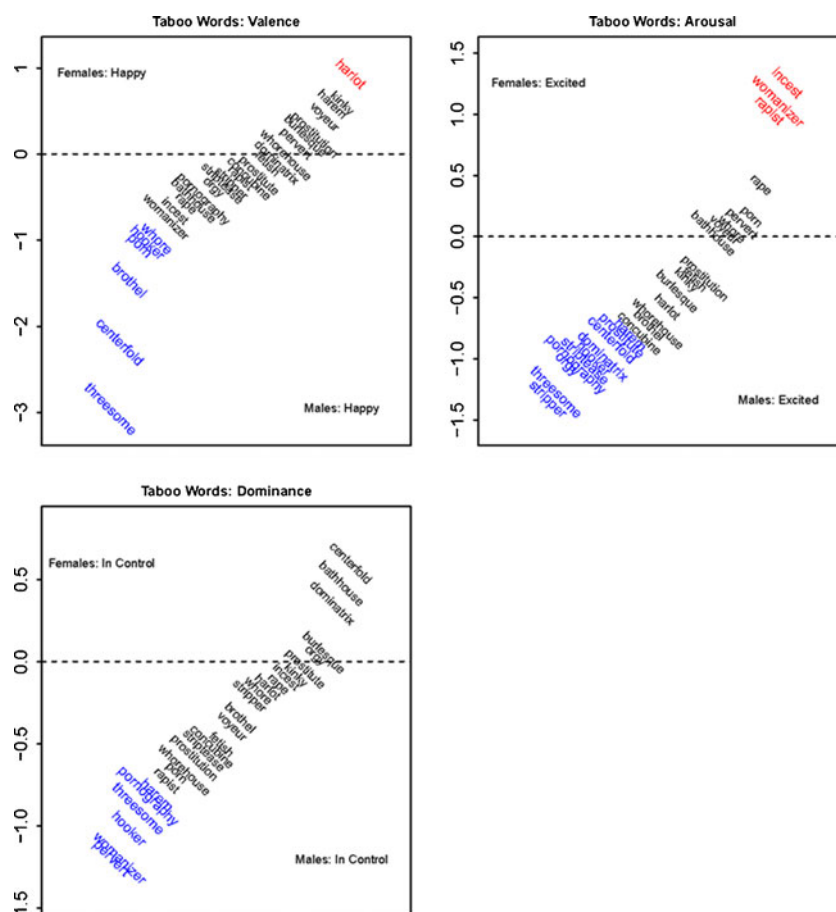


Fig. 11 Gender differences in ratings for taboo words



due to reduced statistical power ($df = 22$). Some interesting contrasts can be seen that might prove interesting to social scientists. For example, both the words “police officer” and “firefighter” are rated as highly arousing, but “police officer” is viewed negatively while “firefighter” is viewed positively. In contrast, “librarian” is a positive but completely unarousing occupation term.

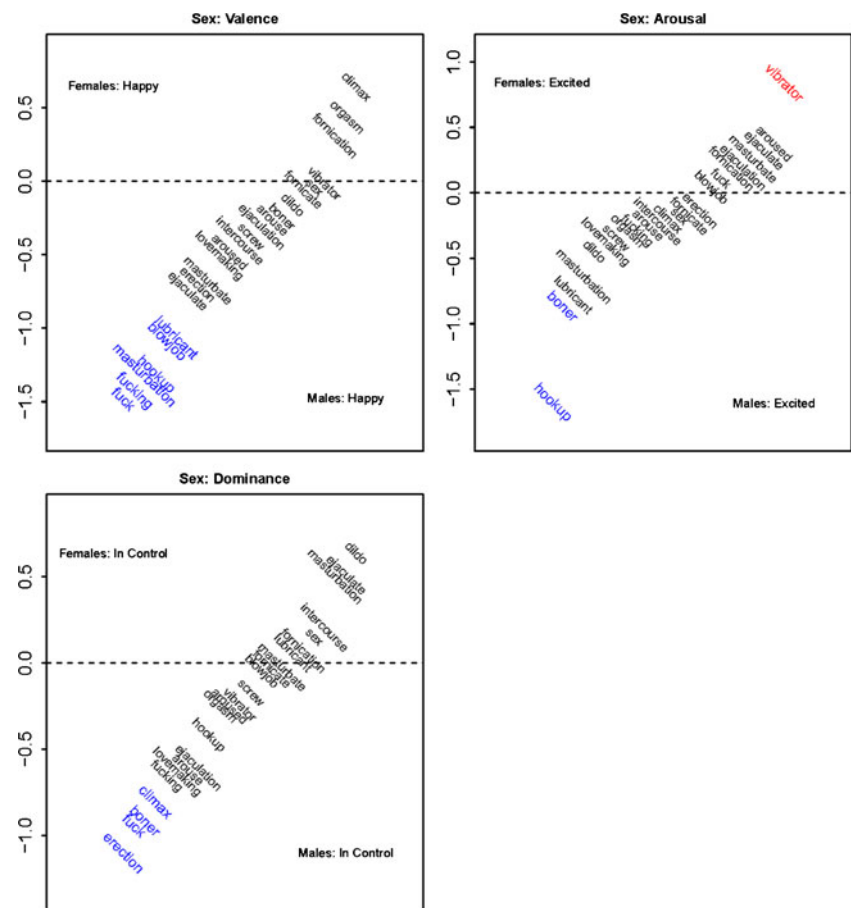
Emotional ratings are also a useful tool for studying gender differences in attitudes and beliefs. Figure 10 reports gender differences in ratings to terms denoting weaponry, with the difference between the ratings of female and male responders on the y-axis. The upper parts of the plots in Fig. 10 show words that were given higher valence, arousal, or dominance ratings by female responders; dotted lines represent the no-difference line. Words in blue color stand for items for which the difference in ratings between gender groups reached significance at the $p < .01$ level in two-tailed independent t tests.

All three emotional dimensions showed a significantly greater number of ratings in the lower parts of the plots (all p values in chi-square tests $< .01$). This indicates that male responders generally have a happier, more aroused, and more in-control attitude toward weapons, especially fire weapons and the bow, for which the gender difference in ratings reached significance.

A similar bias toward higher valence, arousal, and dominance can be observed in ratings of male responders to taboo words and sexual terms. As Figs. 11 and 12 demonstrate, most lexical items in this subset are located below the dotted lines, revealing overall higher ratings for taboo words in male responders (marked in blue if reaching significance) and, in rare cases, in female responders (marked in red if reaching significance). The observed discrepancies in attitudes are corroborated by Janschewitz (2008), Newman, Groom, Handelman, and Pennebaker (2008), and Petersen and Hyde (2010). The discrepancies also explain the disproportionate presence of sexual terms and taboo words among lexical items with exceedingly variable ratings (see the highlighted words in Fig. 2 whose standard deviations are larger than the value predicted from their means).

General discussion

Technological advances are rapidly changing the tools that language researchers have at their disposal. Two main, complementary developments are (1) the collection of large sets of human data through crowdsourcing platforms and (2) the automatic calculation of word characteristics on the basis of

Fig. 12 Gender differences in ratings for sex-related words

relationships between words. In the former case, the current means of digital communication can be used to reach a large audience at an affordable price. The present study is a typical example of this: Instead of having to limit the list of words to a few hundred, because of a lack of human respondents, we extended the list to nearly 14,000 (see Kuperman et al., 2012, for another example of a large-sample rating obtained via crowdsourcing). Our collection of primary demographic information, such as age, gender, and education, additionally enabled refined analyses of both the central tendency and variability in each of the emotional dimensions. Likewise, it paved the way for characterization of attitudes and opinions in the population at large, as well as in specific groups of respondents.

The derivation of word features by means of counting word co-occurrences is an approach that is likely to expand considerably in the coming years. Arguably, the showcase at the moment is the derivation of word meanings by establishing which words co-occur in texts and bits of discourse. Estimates based on word co-occurrences correlate reasonably well with human-generated word associations and semantic similarity ratings. This approach was initiated by Landauer and Dumais (1997) and Burgess (1998). Recent reviews and extensions can be found in Shaoul and Westbury (2010) and Zhao, Li, and Kohonen (2011). The enterprise critically depends on

algorithms that automatically extract word information from collections of texts and calculate various measures of co-occurrence.

Bestgen and Vincze (2012) applied this approach to the affective dimensions of words. They calculated affective norms for over 17,000 words by comparing each word to the thousand words from the ANEW list. The score of each word was derived from the ANEW norms of the words with the closest distance in semantic space. Bestgen and Vincze observed that performance was best when the 30 closest neighbors of the target word were used. This led to correlations of $r = .71$ between the automatically derived values of valence and the human ratings, $r = .56$ for arousal, and $r = .60$ for dominance. All things being equal, these correlations depend on the number of so-called “seed words”—words with known values to which the new words can be compared. The more seed words, the better the estimates for the remaining words. On the other hand, the more seed words for which human data are available, the less the need for automatic extraction of such information. Our extensive data set clearly contributes to the accuracy of such computational estimates. Additionally, it introduces the opportunity to make estimates of textual sentiment for specific reader profiles: for instance, low-educated men, older women, or highly educated youngsters. This in turn may inform the creation of texts that are made

more or less emotionally appealing or arousing to specific target populations.

To sum up, our collection of emotion norms for nearly 14,000 words gives computational and experimental researchers of language use a much wider selection of materials for their studies. Depending on the size of a person's vocabulary, our sample size is estimated to be between one half and one quarter of the words known to individuals. Reliable ratings of the affective states invoked by this number of words will advance the study of the interplay between language and emotion.

Availability

Our ratings are available as supplementary materials for this article and are provided in .csv format. Every value is reported three times: once for each dimension, prefixed with V for valence, A for arousal, and D for dominance. For each word, we report the overall mean (Mean.Sum), standard deviation (SD.Sum), and number of contributing ratings (Rat.Sum). We also report these values for group differences, replacing the suffix .Sum with the following suffixes: .M = male; .F = female; .O = older; .Y = younger; .H = high education; .L = low education. Words are presented in alphabetical order.

We note that group differences (gender, education level, and age), while interesting, are actually quite limited. Taking a conservative $p < .01$ as our definition of a significant difference, fewer than 100 words per dimension meet this criterion (education and arousal include more, with nearly 200 words each). In terms of gender, the differences seem to occur primarily in categories related to sex, violence, and other taboo topics. When these stereotypical domains are under investigation, we do advise people to consider gender differences in the ratings. The semantic categories for other group differences were more difficult to define. In general, unless there is an already established reason to consider group differences, using the overall Sum ratings is, we feel, completely valid.

Author Note This study was supported by an Odysseus Grant awarded by the Government of Flanders (the Dutch-speaking, northern half of Belgium) to M.B., and by an Insight Development Grant from the Social Sciences and Humanities Research Council and a Discovery Grant of the Natural Sciences and Engineering Research Council of Canada to V.K. We thank Dan Wright and Barbara Juhasz for insightful comments on the previous draft.

References

- Altarriba, J., Bauer, L. M., & Benvenuto, C. (1999). Concreteness, context availability, and imageability ratings and word associations for abstract, concrete, and emotion words. *Behavior Research Methods*, 31, 578–602. doi:10.3758/BF03200738
- Amsel, B. D., Urbach, T. P., & Kutas, M. (2012). Perceptual and motor attribute ratings for 559 object concepts. *Behavior Research Methods*, 44, 1028–1041. doi:10.3758/s13428-012-0215-z
- Augustine, A. A., Mehl, M. R., & Larsen, R. J. (2011). A positivity bias in written and spoken English and its moderation by personality and gender. *Social Psychological and Personality Science*, 2, 508–515.
- Baayen, R. H., Feldman, L. B., & Schreuder, R. (2006). Morphological influences on the recognition of monosyllabic monomorphemic words. *Journal of Memory and Language*, 55, 290–313. doi:10.1016/j.jml.2006.03.008
- Balota, D. A., Yap, M. J., Cortese, M. J., Hutchison, K. A., Kessler, B., Loftis, B., et al. (2007). The English lexicon project. *Behavior Research Methods*, 39, 445–459. doi:10.3758/BF03193014
- Bestgen, Y., & Vincze, N. (2012). Checking and bootstrapping lexical norms by means of word similarity indexes. *Behavior Research Methods*, 44, 998–1006. doi:10.3758/s13428-012-0195-z
- Bradley, M. M., & Lang, P. J. (1999). *Affective norms for English words (ANEW): Stimuli, instruction manual and affective ratings (Technical Report No. C-1)*. Gainesville, FL: University of Florida, NIMH Center for Research in Psychophysiology.
- Brems, C., Johnson, M. E., Warner, T. D., & Roberts, L. W. (2010). Health care providers' reports of perceived stigma associated with HIV and AIDS in rural and urban communities. *Journal of HIV/AIDS and Social Services*, 9, 356–370.
- Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis: A critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior Research Methods*, 41, 977–990. doi:10.3758/BRM.41.4.977
- Brysbaert, M., New, B., & Keuleers, E. (2012). Adding part-of-speech information to the SUBTLEX-US word frequencies. *Behavior Research Methods*, 44, 991–997. doi:10.3758/s13428-012-0190-4
- Bureau of Labor Statistics. (2011, May). *National occupational employment and wage estimates: United States*. Retrieved August 31, 2012, from www.bls.gov/oes/current/oes_nat.htm#00-0000
- Burgess, C. (1998). From simple associations to the building blocks of language: Modeling meaning in memory with the HAL model. *Behavior Research Methods, Instruments, & Computers*, 30, 188–198. doi:10.3758/BF03200643
- Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112, 155–159. doi:10.1037/0033-2909.112.1.155
- Cortese, M. J., & Fugett, A. (2004). Imageability ratings for 3,000 monosyllabic words. *Behavior Research Methods, Instruments, & Computers*, 36, 384–387. doi:10.3758/BF03195585
- Eilola, T. M., & Havelka, J. (2010). Affective norms for 210 British English and Finnish nouns. *Behavior Research Methods*, 42, 134–140. doi:10.3758/BRM.42.1.134
- Ferrand, L., New, B., Brysbaert, M., Keuleers, E., Bonin, P., Méot, A., et al. (2010). The French lexicon project: Lexical decision data for 38,840 French words and 38,840 pseudowords. *Behavior Research Methods*, 42, 488–496. doi:10.3758/BRM.42.2.488
- Fraga, I., Piñeiro, A., Acuña-Fariña, C., Redondo, J., & García-Orza, J. (2012). Emotional nouns affect attachment decisions in sentence completion tasks. *Quarterly Journal of Experimental Psychology*, 65, 1740–1759. doi:10.1080/17470218.2012.662989
- Garcia, D., Garas, A., & Schweitzer, F. (2012). Positive words carry less information than negative words. *EPJ Data Science*, 1. doi:10.1140/epjds3
- Gilhooly, K. J., & Logie, R. H. (1980). Age-of-acquisition, imagery, concreteness, familiarity, and ambiguity measures for 1,944 words. *Behavior Research Methods & Instrumentation*, 12, 395–427. doi:10.3758/BF03201693
- Janschewitz, K. (2008). Taboo, emotionally valenced, and emotionally neutral word norms. *Behavior Research Methods*, 40, 1065–1074. doi:10.3758/BRM.40.4.1065

- Juhasz, B. J., & Yap, M. J. (2012). Sensory experience ratings for over 5,000 mono- and disyllabic words. *Behavior Research Methods*. doi:10.3758/s13428-012-0242-9
- Juhasz, B. J., Yap, M. J., Dicke, J., Taylor, S. C., & Gullick, M. M. (2011). Tangible words are recognized faster: The grounding of meaning in sensory and perceptual systems. *Quarterly Journal of Experimental Psychology*, 64, 1683–1691. doi:10.1080/17470218.2011.605150
- Keuleers, E., Brysbaert, M., & New, B. (2010). SUBTLEX-NL: A new measure for Dutch word frequency based on film subtitles. *Behavior Research Methods*, 42, 643–650. doi:10.3758/BRM.42.3.643
- Keuleers, E., Lacey, P., Rastle, K., & Brysbaert, M. (2012). The British lexicon project: Lexical decision data for 28,730 monosyllabic and disyllabic English words. *Behavior Research Methods*, 44, 287–304. doi:10.3758/s13428-011-0118-4
- Kloumann, I. M., Danforth, C. M., Harris, K. D., Bliss, C. A., & Dodds, P. S. (2012). Positivity of the English language. *PLoS One*, 7, e29484. doi:10.1371/journal.pone.0029484
- Kousta, S. T., Vigliocco, G., Vinson, D. P., Andrews, M., & Del Campo, E. (2011). The representation of abstract words: Why emotion matters. *Journal of Experimental Psychology: General*, 140, 14–34. doi:10.1037/a0021446
- Kousta, S. T., Vinson, D. P., & Vigliocco, G. (2009). Emotion words, regardless of polarity, have a processing advantage over neutral words. *Cognition*, 112, 473–481. doi:10.1016/j.cognition.2009.06.007
- Kuperman, V., Stadthagen-Gonzalez, H., & Brysbaert, M. (2012). Age-of-acquisition ratings for 30,000 English words. *Behavior Research Methods*, 44, 978–990. doi:10.3758/s13428-012-0210-4
- 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. doi:10.1037/0033-295X.104.2.211
- Leveau, N., Jhean-Larose, S., Denhière, G., & Nguyen, B. L. (2012). Validating an interlingual metanorm for emotional analysis of texts. *Behavior Research Methods*, 44, 1007–1014. doi:10.3758/s13428-012-0208-y
- Liu, B. (2012). *Sentiment analysis and opinion mining*. San Rafael, CA: Morgan & Claypool.
- Medler, D., Arnoldussen, A., Binder, J., & Seidenberg, M. (2005). *The Wisconsin perceptual attribute ratings database*. Retrieved from www.neuro.mcw.edu/ratings/
- MetLife Foundation. (2011). *What America thinks: MetLife Foundation Alzheimer's survey*. Retrieved August 31, 2012, from www.metlife.com/assets/cao/contributions/foundation/alzheimers-2011.pdf
- Mohammad, S. M., & Turney, P. D. (2010). Emotions evoked by common words and phrases: Using Mechanical Turk to create an emotion lexicon. In *Proceedings of the NAACL-HLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text* (pp. 26–34). New York, NY: Association for Computational Linguistics.
- Moors, A., De Houwer, J., Hermans, D., Wanmaker, S., van Schie, K., Van Harmelen, A. L., . . . Brysbaert, M. (2012). Norms of valence, arousal, dominance, and age of acquisition for 4,300 Dutch words. *Behavior Research Methods*. doi:10.3758/s13428-012-0243-8
- Newman, M. L., Groom, C. J., Handelman, L. D., & Pennebaker, J. W. (2008). Gender differences in language use: An analysis of 14,000 text samples. *Discourse Processes*, 45, 211–236. doi:10.1080/01638530802073712
- Osgood, C. E., Suci, G. J., & Tannenbaum, P. H. (1957). *The measurement of meaning*. Urbana, IL: University of Illinois Press.
- Petersen, J. L., & Hyde, J. S. (2010). A meta-analytic review of research on gender differences in sexuality, 1993–2007. *Psychological Bulletin*, 136, 21–38. doi:10.1037/a0017504
- Rammstedt, B., & Krebs, D. (2007). Does response scale format affect the answering of personality scales? *European Journal of Psychological Assessment*, 23, 32–38.
- Redondo, J., Fraga, I., Padrón, I., & Comesaña, M. (2007). The Spanish adaptation of ANEW (Affective Norms for English Words). *Behavior Research Methods*, 39, 600–605. doi:10.3758/BF03193031
- Schock, J., Cortese, M. J., & Khanna, M. M. (2012). Imageability estimates for 3,000 disyllabic words. *Behavior Research Methods*, 44, 374–379. doi:10.3758/s13428-011-0162-0
- Scott, G. G., O'Donnell, P. J., & Sereno, S. C. (2012). Emotion words affect eye fixations during reading. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 38, 783–792. doi:10.1037/a0027209
- Shaoul, C., & Westbury, C. (2010). Exploring lexical co-occurrence space using HiDex. *Behavior Research Methods*, 42, 393–413. doi:10.3758/BRM.42.2.393
- Soares, A. P., Comesaña, M., Pinheiro, A. P., Simões, A., & Frade, C. S. (2012). The adaptation of the Affective Norms for English Words (ANEW) for European Portuguese. *Behavior Research Methods*, 44, 256–269. doi:10.3758/s13428-011-0131-7
- Stadthagen-Gonzalez, H., & Davis, C. J. (2006). The Bristol norms for age of acquisition, imageability, and familiarity. *Behavior Research Methods*, 38, 598–605. doi:10.3758/BF03193891
- Tillotson, S. M., Siakaluk, P. D., & Pexman, P. M. (2008). Body-object interaction ratings for 1,618 monosyllabic nouns. *Behavior Research Methods*, 40, 1075–1078. doi:10.3758/BRM.40.4.1075
- Van Overschelde, J. P., Rawson, K. A., & Dunlosky, J. (2004). Category norms: An updated and expanded version of the Battig and Montague (1969) norms. *Journal of Memory and Language*, 50, 289–335. doi:10.1016/j.jml.2003.10.003
- Verona, E., Sprague, J., & Sadeh, N. (2012). Inhibitory control and negative emotional processing in psychopathy and antisocial personality disorder. *Journal of Abnormal Psychology*, 121, 498–510.
- YouGov. (2011). *Cancer Britons most feared disease*. Retrieved August 31, 2012, from <http://yougov.co.uk/news/2011/08/15/cancer-britons-most-feared-disease/>
- Zhao, X., Li, P., & Kohonen, T. (2011). Contextual self-organizing map: Software for constructing semantic representations. *Behavior Research Methods*, 43, 77–88. doi:10.3758/s13428-010-0042-z