

Word embeddings quantify 100 years of gender and ethnic stereotypes

Author(s): Nikhil Garg, Londa Schiebinger, Dan Jurafsky and James Zou

Source: *Proceedings of the National Academy of Sciences of the United States of America*, April 17, 2018, Vol. 115, No. 16 (April 17, 2018), pp. E3635-E3644

Published by: National Academy of Sciences

Stable URL: <https://www.jstor.org/stable/10.2307/26508864>

REFERENCES

Linked references are available on JSTOR for this article:

https://www.jstor.org/stable/10.2307/26508864?seq=1&cid=pdf-reference#references_tab_contents

You may need to log in to JSTOR to access the linked references.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



JSTOR

National Academy of Sciences is collaborating with JSTOR to digitize, preserve and extend access to *Proceedings of the National Academy of Sciences of the United States of America*



Word embeddings quantify 100 years of gender and ethnic stereotypes

Nikhil Garg^{a,1}, Londa Schiebinger^b, Dan Jurafsky^{c,d}, and James Zou^{e,f,1}

^aDepartment of Electrical Engineering, Stanford University, Stanford, CA 94305; ^bDepartment of History, Stanford University, Stanford, CA 94305;

^cDepartment of Linguistics, Stanford University, Stanford, CA 94305; ^dDepartment of Computer Science, Stanford University, Stanford, CA 94305;

^eDepartment of Biomedical Data Science, Stanford University, Stanford, CA 94305; and ^fChan Zuckerberg Biohub, San Francisco, CA 94158

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved March 12, 2018 (received for review November 22, 2017)

Word embeddings are a powerful machine-learning framework that represents each English word by a vector. The geometric relationship between these vectors captures meaningful semantic relationships between the corresponding words. In this paper, we develop a framework to demonstrate how the temporal dynamics of the embedding helps to quantify changes in stereotypes and attitudes toward women and ethnic minorities in the 20th and 21st centuries in the United States. We integrate word embeddings trained on 100 y of text data with the US Census to show that changes in the embedding track closely with demographic and occupation shifts over time. The embedding captures societal shifts—e.g., the women's movement in the 1960s and Asian immigration into the United States—and also illuminates how specific adjectives and occupations became more closely associated with certain populations over time. Our framework for temporal analysis of word embedding opens up a fruitful intersection between machine learning and quantitative social science.

word embedding | gender stereotypes | ethnic stereotypes

The study of gender and ethnic stereotypes is an important topic across many disciplines. Language analysis is a standard tool used to discover, understand, and demonstrate such stereotypes (1–5). Previous literature broadly establishes that language both reflects and perpetuates cultural stereotypes. However, such studies primarily leverage human surveys (6–16), dictionary and qualitative analysis (17), or in-depth knowledge of different languages (18). These methods often require time-consuming and expensive manual analysis and may not easily scale across types of stereotypes, time periods, and languages. In this paper, we propose using word embeddings, a commonly used tool in natural language processing (NLP) and machine learning, as a framework to measure, quantify, and compare beliefs over time. As a specific case study, we apply this tool to study the temporal dynamics of gender and ethnic stereotypes in the 20th and 21st centuries in the United States.

In word-embedding models, each word in a given language is assigned to a high-dimensional vector such that the geometry of the vectors captures semantic relations between the words—e.g., vectors being closer together has been shown to correspond to more similar words (19). These models are typically trained automatically on large corpora of text, such as collections of Google News articles or Wikipedia, and are known to capture relationships not found through simple co-occurrence analysis. For example, the vector for France is close to vectors for Austria and Italy, and the vector for Xbox is close to that of PlayStation (19). Beyond nearby neighbors, embeddings can also capture more global relationships between words. The difference between London and England—obtained by simply subtracting these two vectors—is parallel to the vector difference between Paris and France. This pattern allows embeddings to capture analogy relationships, such as London to England is as Paris to France.

Recent works demonstrate that word embeddings, among other methods in machine learning, capture common stereotypes because these stereotypes are likely to be present, even if subtly,

in the large corpora of training texts (20–23). For example, the vector for the adjective honorable would be close to the vector for man, whereas the vector for submissive would be closer to woman. These stereotypes are automatically learned by the embedding algorithm and could be problematic if the embedding is then used for sensitive applications such as search rankings, product recommendations, or translations. An important direction of research is to develop algorithms to debias the word embeddings (20).

In this paper, we take another approach. We use the word embeddings as a quantitative lens through which to study historical trends—specifically trends in the gender and ethnic stereotypes in the 20th and 21st centuries in the United States. We develop a systematic framework and metrics to analyze word embeddings trained over 100 y of text corpora. We show that temporal dynamics of the word embedding capture changes in gender and ethnic stereotypes over time. In particular, we quantify how specific biases decrease over time while other stereotypes increase. Moreover, dynamics of the embedding strongly correlate with quantifiable changes in US society, such as demographic and occupation shifts. For example, major transitions in the word embedding geometry reveal changes in the descriptions of genders and ethnic groups during the women's movement in the 1960s–1970s and Asian-American population growth in the 1960s and 1980s. We validate our findings on external metrics and show that our results are robust to the different algorithms for training the word embeddings. Our framework reveals and quantifies how stereotypes toward women and ethnic groups have evolved in the United States.

Significance

Word embeddings are a popular machine-learning method that represents each English word by a vector, such that the geometry between these vectors captures semantic relations between the corresponding words. We demonstrate that word embeddings can be used as a powerful tool to quantify historical trends and social change. As specific applications, we develop metrics based on word embeddings to characterize how gender stereotypes and attitudes toward ethnic minorities in the United States evolved during the 20th and 21st centuries starting from 1910. Our framework opens up a fruitful intersection between machine learning and quantitative social science.

Author contributions: N.G., L.S., D.J., and J.Z. designed research; N.G. and J.Z. performed research; and N.G. and J.Z. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Published under the PNAS license.

Data deposition: Data and code related to this paper are available on GitHub (<https://github.com/nikhgarg/EmbeddingDynamicStereotypes>).

¹To whom correspondence may be addressed. Email: nkgarg@stanford.edu or jamesz@stanford.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1720347115/-DCSupplemental.

Published online April 3, 2018.

Our results demonstrate that word embeddings are a powerful lens through which we can systematically quantify common stereotypes and other historical trends. Embeddings thus provide an important quantitative metric which complements existing, more qualitative, linguistic and sociological analyses of biases. In *Embedding Framework Overview and Validations*, we validate that embeddings accurately capture sociological trends by comparing associations in the embeddings with census and other externally verifiable data. In *Quantifying Gender Stereotypes* and *Quantifying Ethnic Stereotypes* we apply the framework to quantify the change in stereotypes of women, men, and ethnic minorities. We further discuss our findings in *Discussion* and provide additional details in *Materials and Methods*.

Embedding Framework Overview and Validations

In this section, we briefly describe our methods and data and then validate our findings. We focus on showing that word embeddings are an effective tool to study historical biases and stereotypes by relating measurements from these embeddings to historical census and survey data. The consistent replication of such historical data, both in magnitude and in direction of biases, validates the use of embeddings in such work. This section extends the analysis of refs. 20 and 21 in showing that embeddings can also be used as a comparative tool over time as a consistent metric for various biases.

Summary of Data and Methods. We now briefly describe our datasets and methods, leaving details to *Materials and Methods* and *SI Appendix, section A*. All of our code and embeddings are available publicly*. For contemporary snapshot analysis, we use the standard Google News word2vec vectors trained on the Google News dataset (24, 25). For historical temporal analysis, we use previously trained Google Books/Corpus of Historical American English (COHA) embeddings, which are a set of nine embeddings, each trained on a decade in the 1900s, using the COHA and Google Books (26). As additional validation, we train, using the GloVe algorithm (27), embeddings from the *New York Times* Annotated Corpus (28) for every year between 1988 and 2005. We then collate several word lists to represent each gender[†] (men, women) and ethnicity[‡] (White, Asian, and Hispanic), as well as neutral words (adjectives and occupations). For occupations, we use historical US census data (29) to extract the percentage of workers in each occupation that belong to each gender or ethnic group and compare it to the bias in the embeddings.

Using the embeddings and word lists, one can measure the strength of association (embedding bias) between neutral words and a group. As an example, we overview the steps we use to quantify the occupational embedding bias for women. We first compute the average embedding distance between words that represent women—e.g., she, female—and words for occupations—e.g., teacher, lawyer. For comparison, we also compute the average embedding distance between words that represent men and the same occupation words. A natural metric for the embedding bias

is the average distance for women minus the average distance for men. If this value is negative, then the embedding more closely associates the occupations with men. More generally, we compute the representative group vector by taking the average of the vectors for each word in the given gender/ethnicity group. Then we compute the average Euclidean distance between each representative group vector and each vector in the neutral word list of interest, which could be occupations or adjectives. The difference of the average distances is our metric for bias—we call this the relative norm difference or simply embedding bias.

We use ordinary least-squares regressions to measure associations in our analysis. In this paper, we report r^2 and the coefficient P value for each regression, along with the intercept confidence interval when relevant.

Validation of the Embedding Bias. To verify that the bias in the embedding accurately reflects sociological trends, we compare the trends in the embeddings with quantifiable demographic trends in the occupation participation, as well as historical surveys of stereotypes. First, we use women and minority ethnic participation statistics (relative to men and Whites, respectively) in different occupations as a benchmark because it is an objective metric of social changes. We show that the embedding accurately captures both gender and ethnic occupation percentages and consistently reflects historical changes.

Next, we validate that the embeddings capture personality trait stereotypes. A difficulty in social science is the relative dearth of historical data to systematically quantify stereotypes, which highlights the value of our embedding framework as a quantitative tool but also makes it challenging to directly confirm our findings on adjectives. Nevertheless, we make use of the best available data from historical surveys, gender stereotypes from 1977 and 1990 (6, 7) and ethnic stereotypes from the Princeton trilogy from 1933, 1951, and 1969 (8–10).

Comparison with women's occupation participation. We investigate how the gender bias of occupations in the word embeddings relates to the empirical percentage of women in each of these occupations in the United States. Fig. 1 shows, for each occupation, the relationship between the relative percentage (of women) in the occupation in 2015 and the relative norm distance between words associated with women and men in the Google News embeddings. (Occupations whose 2015 percentage is not available, such as midwife, are omitted. We further note that the Google News embedding is trained on a corpus



Fig. 1. Women's occupation relative percentage vs. embedding bias in Google News vectors. More positive indicates more associated with women on both axes. $P < 10^{-10}$, $r^2 = 0.499$. The shaded region is the 95% bootstrapped confidence interval of the regression line. In this single embedding, then, the association in the embedding effectively captures the percentage of women in an occupation.

* All of our own data and analysis tools are available on GitHub at <https://github.com/nikhgarg/EmbeddingDynamicStereotypes>. Census data are available through the Integrated Public Use Microdata Series (29). We link to the sources for each embedding used in *Materials and Methods*.

[†] There is an increasingly recognized difference between sex and gender and thus between the words male/female and man/woman, as well as nonbinary categories. We limit our analysis to the two major binary categories due to technical limitations, and we use male and female as part of the lists of words associated with men and women, respectively, when measuring gender associations. We also use results from refs. 6 and 7 which study stereotypes associated with sex.

[‡] When we refer to Whites or Asians, we specifically mean the non-Hispanic subpopulation. For each ethnicity, we generate a list of common last names among the group. Unfortunately, our present methods do not extend to Blacks due to large overlaps in common last names among Whites and Blacks in the United States.

interval $(-0.027, -0.001)$ are significantly correlated with the embedding bias.

Next, we conduct two additional separate regressions to test that the embedding bias captures the same extra stereotype information as do the crowdscore scores, information that is missing in the census data. In each regression, the occupation percentage difference is the independent covariate. In one, the embedding bias is the dependent variable; in the other, stereotype score is. In these regressions, a negative (positive) residual indicates that the embedding bias or stereotype score is closer to words associated with women (men) than is to be expected given the gender percentages in the occupation. We find that the residuals between the two regressions correlate significantly (Pearson coefficient 0.811, $P < 10^{-10}$). This correlation suggests that the embedding bias captures the crowdscore human stereotypes beyond that which can be explained by empirical differences in occupation proportions.

Where such crowdsourcing is not possible, such as in studying historical biases, word embeddings can thus further serve as an effective measurement tool. Further, although the analysis in the previous section shows a strong relationship between census data and embedding bias, it is important to note that biases beyond census data also appear in the embedding.

Quantifying Changing Attitudes with Adjective Embeddings. We now apply the insight that embeddings can be used to make comparative statements over time to study how the description of women—through adjectives—in literature and the broader culture has changed over time. Using word embeddings to analyze biases in adjectives could be an especially useful approach because the literature is lacking systematic and quantitative metrics for adjective biases. We find that—as a whole—portrayals have changed dramatically over time, including for the better in some measurable ways. Furthermore, we find evidence for how the women's movement in the 1960s and 1970s led to a systemic change in such portrayals.

How overall portrayals change over time. We first establish that comparing the embeddings over time could reveal global shifts in society in regard to gender portrayals. Fig. 4 shows the Pearson correlation in embedding bias scores for adjectives over time between COHA embeddings for each pair of decades. As expected, the highest correlation values are near the diagonals; embeddings (and attitudes) are most similar to those from adjacent decades. More strikingly, the matrix exhibits two clear blocks. There is a sharp divide between the 1960s and 1970s, the height of the women's movement in the United States, during which there was a large push to end legal and social barriers for women in education and the workplace (32, 33). The transition in the gender embeddings from 1960 to 1970 is statistically significant

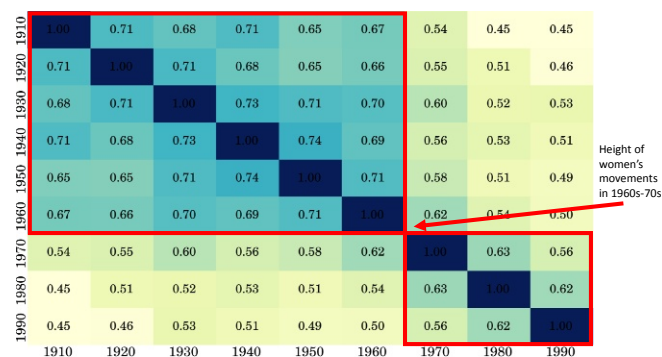


Fig. 4. Pearson correlation in embedding bias scores for adjectives over time between embeddings for each decade. The phase shift in the 1960s–1970s corresponds to the US women's movement.

icant ($P < 10^{-4}$, Kolmogorov–Smirnov two-sample test) and is larger than the change between any two other adjacent decades. See *SI Appendix, section B.3.3* for a more detailed description of the test and all statistics.

We note that the effects of the women's movement, including on inclusive language, are well documented (18, 33–36); this work provides a quantitative way to measure the rate and extent of the change. A potential extension and application of this work would be to study how various narratives and descriptions of women developed and competed over time.

Individual words whose biases changed over time. As an example of such work, we consider a subset of the adjectives describing competence, such as intelligent, logical, and thoughtful (see *SI Appendix, section A.3* for a full list of words; these words were curated from various online sources). Since the 1960s, this group of words on average has increased in association with women over time (from strongly biased toward men to less so): In a regression with embedding bias from each word as the dependent variable and years from 1960 to 1990 as the covariate, the coefficient is positive; i.e., there is a (small) positive trend (0.005 increase in women association per decade, $P = 0.0036$). At this rate, such adjectives would be equally associated with women as with men a little after the year 2020.

As a comparison, we also analyze a subset of adjectives describing physical appearance—e.g., attractive, ugly, and fashionable—and the bias of these words did not change significantly since the 1960s (null hypothesis of no trend not rejected with $P > 0.2$). Although the trend regarding intelligence is encouraging, the top adjectives are still potentially problematic, as displayed in Table 2.

We note that this analysis is an exploration; perceived competence and physical appearance are just two components of gender stereotypes. Models in the literature suggest that stereotypes form along several dimensions, e.g., warmth and competence (16). A more complete analysis would first collect externally validated lists of words that describe each such dimension and then measure the embedding association with respect to these lists over time.

The embedding also reveals interesting patterns in how individual words evolve over time in their gender association. For example, the word hysterical used to be, until the mid-1900s, a catchall term for diagnosing mental illness in women but has since become a more general word (37); such changes are clearly reflected in the embeddings, as hysterical fell from a top 5 woman-biased word in 1920 to not in the top 100 in 1990 in the COHA embeddings[#]. On the other hand, the word emotional becomes much more strongly associated with women over time in the embeddings, reflecting its current status as a word that is largely associated with women in a pejorative sense (38).

These results together demonstrate the value and potential of leveraging embeddings to study biases over time. The embeddings capture subtle individual changes in association, as well as larger historical changes. Overall, they paint a picture of a society with decreasing but still significant gender biases.

Quantifying Ethnic Stereotypes

We now turn our attention to studying ethnic biases over time. In particular we show how immigration and other 20th-century trends broadly influenced how Asians were viewed in the United States. We also show that embeddings can serve as effective tools to analyze finer-grained trends by analyzing the portrayal of Islam in the *New York Times* from 1988 to 2005 in the context of terrorism.

[#]We caution that due to the noisy nature of word embeddings, dwelling on individual word rankings in isolation is potentially problematic. For example, hysterical is more highly associated with women in the Google News vectors than emotional. For this reason we focus on large shifts between embeddings.

Table 2. Top adjectives associated with women in 1910, 1950, and 1990 by relative norm difference in the COHA embedding

1910	1950	1990
Charming	Delicate	Maternal
Placid	Sweet	Morbid
Delicate	Charming	Artificial
Passionate	Transparent	Physical
Sweet	Placid	Caring
Dreamy	Childish	Emotional
Indulgent	Soft	Protective
Playful	Colorless	Attractive
Mellow	Tasteless	Soft
Sentimental	Agreeable	Tidy

Trends in Asian Stereotypes. To study Asian stereotypes in the embeddings, we use common and distinctly Asian last names, identified through a process described in *SI Appendix, section A.2*. This process results in a list of 20 last names that are primarily but not exclusively Chinese last names.

The embeddings illustrate a dramatic story of how Asian-American stereotypes developed and changed in the 20th century. Fig. 5 shows the Pearson correlation coefficient of adjective biases for each pair of embeddings over time. As with gender, the analysis shows how external events changed attitudes. There are two phase shifts in the correlation: in the 1960s, which coincide with a sharp increase in Asian immigration into the United States due to the passage of the 1965 Immigration and Nationality Act, and in the 1980s, when immigration continued and the second-generation Asian-American population emerged (39). Using the same Kolmogorov–Smirnov test on the correlation differences described in the previous section, the phase shifts between the 1950s–1960s ($P = 0.011$) and 1970s–1980s ($P < 10^{-3}$) are significant, while the rest are not ($P > 0.070$).

We extract the most biased adjectives toward Asians (when compared with Whites) to gain more insights into factors driving these global changes in the embedding. Table 3 shows the most Asian-biased adjectives in 1910, 1950, and 1990. Before 1950, strongly negative words, especially those often used to describe outsiders, are among the words most associated with Asians: barbaric, hateful, monstrous, bizarre, and cruel. However, starting around 1950 and especially by 1980, with a rising Asian population in the United States, these words are largely replaced by words often considered stereotypic (40–42) of Asian Americans today: sensitive, passive, complacent, active, and hearty, for example. See *SI Appendix, Table C.8* for the complete list of the top 10 most Asian-associated words in each decade.

Using our methods regarding trends, we can quantify this change more precisely: Fig. 6 shows the relative strength of the Asian association for words used to describe outsiders over time. As opposed to the adjectives overall, which see two distinct phase shifts in Asian association, the words related to outsiders steadily decrease in Asian association over time—except around World War II—indicating that broader globalization trends led to changing attitudes with regard to such negative portrayals. Overall, the word embeddings exhibit a remarkable change in adjectives and attitudes toward Asian Americans during the 20th century.

Trends in Other Ethnic and Cultural Stereotypes. Similar trends appear in other datasets as well. Fig. 7 shows, in the *New York Times* over two decades, how words related to Islam (vs. those related to Christianity) associate with terrorism-related words. Similar to how we measure occupation-related bias, we create a list of words associated with terrorism, such as terror, bomb, and violence. We then measure how associated these words appear to

be in the text to words representing each religion, such as mosque and church, for Islam and Christianity, respectively. (Full word lists are available in *SI Appendix, section A*.) Throughout the time period in the *New York Times*, Islam is more associated with terrorism than is Christianity. Furthermore, an increase in the association can be seen both after the 1993 World Trade Center bombings and after September 11, 2001. With a more recent dataset and using more news outlets, it would be useful to study how such attitudes have evolved since 2005.

We illustrate how word embeddings capture stereotypes toward other ethnic groups. For example, *SI Appendix, Fig. C.4*, with Russian names, shows a dramatic shift in the 1950s, the start of the Cold War, and a minor shift during the initial years of the Russian Revolution in the 1910s–1920s. Furthermore, *SI Appendix, Fig. C.5*, the correlation over time plot with Hispanic names, serves as an effective control group. It shows more steady changes in the embeddings rather than the sharp transitions found in Asian and Russian associations. This pattern is consistent with the fact that numerous events throughout the 20th century influenced the story of Hispanic immigration into the United States, with no single event playing too large a role (43).

These patterns demonstrate the usefulness of our methods to study ethnic as well as gender bias over time; similar analyses can be performed to examine shifts in the attitudes toward other ethnic groups, especially around significant global events. In particular, it would be interesting to more closely measure dehumanization and “othering” of immigrants and other groups using a suite of linguistic techniques, validating and extending the patterns discovered in this work.

Discussion

In this work, we investigate how the geometry of word embeddings, with respect to gender and ethnic stereotypes, evolves over time and tracks with empirical demographic changes in the United States. We apply our methods to analyze word embeddings trained over 100 y of text data. In particular, we quantify the embedding biases for occupations and adjectives. Using occupations allows us to validate the method when the embedding associations are compared with empirical participation rates for each occupation. We show that both gender and ethnic occupation biases in the embeddings significantly track with the actual occupation participation rates. We also show that adjective associations in the embeddings provide insight into how different groups of people are viewed over time.

As in any empirical work, the robustness of our results depends on the data sources and the metrics we choose to represent bias or association. We choose the relative norm difference metric for its simplicity, although many such metrics are reasonable. Refs. 20 and 21 leverage alternate metrics, for example. Our metric agrees with other possible metrics—both

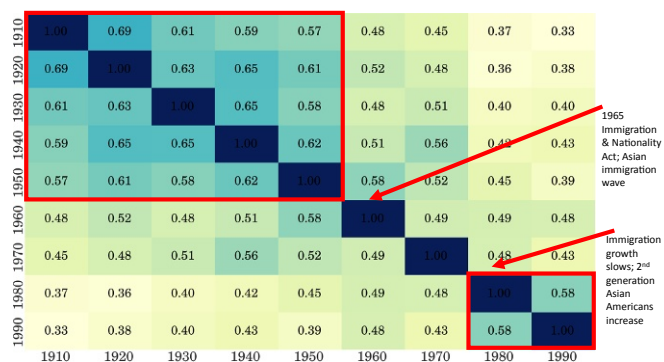


Fig. 5. Pearson correlation in embedding Asian bias scores for adjectives over time between embeddings for each decade.

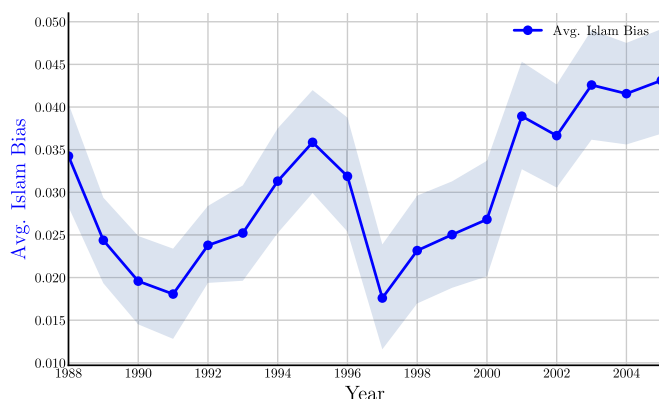


Fig. 7. Religious (Islam vs. Christianity) bias score over time for words related to terrorism in *New York Times* data. Note that embeddings are trained in 3-y windows, so, for example, 2000 contains data from 1999–2001. The shaded region is the bootstrap SE interval.

trained, and between them cover the best-known algorithms to construct embeddings. One finding in this work is that, although there is some heterogeneity, gender and ethnic bias is generally consistent across embeddings. Here we restrict descriptions to embeddings used in the main exposition. For consistency, only single words are used, all vectors are normalized by their l_2 norm, and words are converted to lowercase.

Google News word2vec vectors. Vectors trained on about 100 billion words in the Google News dataset (24, 25). Vectors are available at <https://code.google.com/archive/p/word2vec/>.

Google Books/COHA. Vectors trained on a combined corpus of genre-balanced Google Books and the COHA (48) by the authors of ref. 26. For each decade, a separate embedding is trained from the corpus data corresponding to that decade. The dataset is specifically designed to enable comparisons across decades, and the creators take special care to avoid selection bias issues. The vectors are available at <https://nlp.stanford.edu/projects/histwords/>, and we limit our analysis to the SVD and skip-gram with negative sampling (SGNS) (also known as word2vec) embeddings in the 1900s. Note that the Google Books data may include some non-American sources and the external metrics we use are American. However, this does not appreciably affect results. In the main text, we exclusively use SGNS embeddings; results with SVD embeddings are in *SI Appendix* and are qualitatively similar to the SGNS results. Unless otherwise specified, COHA indicates these embeddings trained using the SGNS algorithm.

New York Times. We train embeddings over time from *The New York Times* Annotated Corpus (28), using 1.8 million articles from the *New York Times* between 1988 and 2005. We use the GloVe algorithm (27) and train embeddings over 3-y windows (so the 2000 embeddings, for example, contain articles from 1999 to 2001).

In *SI Appendix* we also use other embeddings available at <https://nlp.stanford.edu/projects/glove/>.

Related Work. Word embedding was developed as a framework to represent words as a part of the artificial intelligence and natural language processing pipeline (25). Ref. 20 demonstrated that word embeddings capture gender stereotypes, and ref. 21 additionally verified that the embedding accurately reflects human biases by comparing the embedding results with that of the implicit association test. While these two papers analyzed the bias of the static Google News embedding, our paper investigates the temporal changes in word embeddings and studies how embeddings over time capture historical trends. Our paper also studies attitudes toward women and ethnic minorities by quantifying the embedding of adjectives. The focus of ref. 20 is to develop algorithms to reduce the gender stereotype in the embedding, which is important for sensitive applications of embeddings. In contrast, our aim is not to debias, but to leverage the embedding bias to study historical changes that are otherwise challenging to quantify. Ref. 21 shows that embeddings contain each of the associations commonly found in the implicit association test. For example, European-American names are more similar to pleasant (vs. unpleasant) words than are African-American names, and male names are more similar to career (vs. family) words than are female names. Similarly, they show that, in the Google News embeddings, census data correspond to bias in the embeddings for gender.

The study of gender and ethnic stereotypes is a large focus of linguistics and sociology and is too extensive to be surveyed here (1–5). Our main innovation is the use of word embeddings, which provides a unique lens to measure and quantify biases. Another related field in linguistics studies how language changes over time and has also recently used word embeddings as a tool (49–51). However, this literature primarily studies semantic changes, such as how the word gay used to primarily mean cheerful and now means predominantly means homosexual (26, 52), and does not investigate bias.

Word Lists and External Metrics. Two types of word lists are used in this work: group words and neutral words. Group words represent groups of people, such as each gender and ethnicity. Neutral words are those that are not intrinsically gendered or ethnic (for example, fireman or mailman would be gendered occupation titles and so are excluded); relative similarities between neutral words and a pair of groups (such as men vs. women) are used to measure the strength of the association in the embeddings. In this work, we use occupations and various adjective lists as neutral words.

Gender. For gender, we use noun and pronoun pairs (such as he/she, him/her, etc.).

Race/ethnicity. To distinguish various ethnicities, we leverage the fact that the distribution of last names in the United States differs significantly by ethnicity, with the notable exception of White and Black last names. Starting with a breakdown of ethnicity by last name compiled by ref. 53, we identify 20 last names for each ethnicity as detailed in *SI Appendix, section A.2*. Our procedure, however, produces almost identical lists for White and Black Americans (with the names being mostly White by percentage), and so the analysis does not include Black Americans.

Occupation census data. We use occupation words for which we have gender and ethnic subgroup information over time. Group occupation percentages are obtained from the Integrated Public Use Microdata Series (IPUMS), part of the University of Minnesota Historical Census Project (29). Data coding and preprocessing are done as described in ref. 44, which studies wage dynamics as women enter certain occupations over time. The IPUMS dataset includes a column, OCC1950, coding occupation census data as it would have been coded in 1950, allowing accurate interyear analysis. We then hand map the occupations from this column to single-word occupations (e.g., chemical engineer and electrical engineer both become engineer, and chemist is counted as both chemist and scientist) and hand code a subset of the occupations as professional. In all plots containing occupation percentages for gender, we use the percentage difference between women and men in the occupation:

$$p_{\text{women}} - p_{\text{men}}$$

where p_{women} = % of occupation that is women
 p_{men} = % of occupation that is men.

For ethnicity, we similarly report the percentage difference, except we first condition on the workers being in one of the groups in question:

$$\frac{p_{\text{min}} - p_{\text{white}}}{p_{\text{min}} + p_{\text{white}}}$$

where p_{min} = % of occupation that is minority group in question
 p_{white} = % of occupation that is White.

In each case, a value of 0 indicates an equal number of each group in the occupation. We note that the results do not qualitatively change if instead the logit proportion (or conditional logit proportion) of the minority group is used, as in ref. 44 (*SI Appendix, section A.6*).

Occupation gender stereotypes. For a limited set of occupations, we use gender stereotype scores collected from users on Amazon Mechanical Turk by ref. 20. These scores are compared with embedding gender association.

Adjectives. To study associations with adjectives over time, several separate lists are used. To compare gender adjective embedding bias to external metrics, we leverage a list of adjectives labeled by how stereotypically associated with men or women they are, as determined by a group of subjects in 1977 and 1990 (6, 7). For Chinese adjective embedding bias, we use a list of stereotypes from the Princeton trilogy (8–10). For all other analyses using adjectives, a larger list of adjectives is used, primarily from ref. 54. Except when otherwise specified, adjectives are used to refer to this larger list.

Metrics. Given two vectors, their similarity can be measured either by their negative difference norm, as in Eq. 1, or by their cosine similarity, as in Eq. 2. The denominators are omitted because all vectors have norm 1:

$$\text{neg-norm-dif}(u, v) = -\|u - v\|_2 \quad [1]$$

$$\text{cos-sim}(u, v) = u \cdot v. \quad [2]$$

The association between the group words and neutral words is calculated as follows: Construct a group vector by averaging the vectors for each word in the group; then calculate the similarity between this average vector and each word in the neutral list as above.

The relative norm distance, which captures the relative strength of association of a set of neutral words with respect to two groups, is as described in Eq. 3, where M is the set of neutral word vectors, v_1 is the average vector for group one, and v_2 is the average vector for group two. The more positive (negative) that the relative norm distance is, the more associated

the neutral words are toward group two (one). In this work, when we say that a word is biased toward a group with respect to another group, we specifically mean in the context of the relative norm distance. Bias score also refers to this metric:

$$\text{relative norm distance} = \sum_{v_m \in M} \|v_m - v_1\|_2 - \|v_m - v_2\|_2. \quad [3]$$

We can also use cosine similarity rather than the Euclidean 2-norm. [SI Appendix, section A.5](#) shows that the choice of similarity measure is not important; the respective metrics using each similarity measure correlate highly with one another (Pearson coefficient > 0.95 in most cases). In the main text, we exclusively use the relative norm.

ACKNOWLEDGMENTS. J.Z. is supported by a Chan–Zuckerberg Biohub Investigator grant and National Science Foundation (NSF) Grant CRII 1657155. N.G. is supported by the NSF Graduate Research Fellowship under Grant DGE-114747.

- Hamilton DL, Troler TK (1986) *Stereotypes and Stereotyping: An Overview of the Cognitive Approach in Prejudice, Discrimination, and Racism* (Academic, San Diego), pp 127–163.
- Basow SA (1992) *Gender: Stereotypes and Roles* (Thomson Brooks/Cole Publishing Co, Belmont, CA), 3rd Ed.
- Wetherell M, Potter J (1992) *Mapping the Language of Racism: Discourse and the Legitimation of Exploitation* (Columbia Univ Press, New York).
- Holmes J, Meyerhoff M, eds (2004) *The Handbook of Language and Gender* (Blackwell Publishing Ltd, Oxford).
- Coates J (2016) *Women, Men and Language: A Sociolinguistic Account of Gender Differences in Language* (Routledge, London).
- Williams JE, Best DL (1977) Sex stereotypes and trait favorability on the adjective check list. *Educ Psychol Meas* 37:101–110.
- Williams JE, Best DL (1990) *Measuring Sex Stereotypes: A Multination Study* (Sage Publications, Thousand Oaks, CA), Rev Ed.
- Katz D, Braly K (1933) Racial stereotypes of one hundred college students. *J Abnorm Soc Psychol* 28:280–290.
- Gilbert GM (1951) Stereotype persistence and change among college students. *J Abnorm Soc Psychol* 46:245–254.
- Karlins M, Coffman TL, Walters G (1969) On the fading of social stereotypes: Studies in three generations of college students. *J Pers Soc Psychol* 13:1–16.
- Devine PG, Elliot AJ (1995) Are racial stereotypes really fading? The Princeton trilogy revisited. *Pers Soc Psychol Bull* 21:1139–1150.
- Diekmann AB, Eagly AH (2000) Stereotypes as dynamic constructs: Women and men of the past, present, and future. *Pers Soc Psychol Bull* 26:1171–1188.
- Bergsieker HB, Leslie LM, Constantine VS, Fiske ST (2012) Stereotyping by omission: Eliminate the negative, accentuate the positive. *J Pers Soc Psychol* 102:1214–1238.
- Madon S, et al. (2001) Ethnic and national stereotypes: The Princeton trilogy revisited and revised. *Pers Soc Psychol Bull* 27:996–1010.
- Gaertner SL, Dovidio JF (1986) *The Aversive Form of Racism* (Academic, San Diego).
- Fiske ST, Cuddy AJC, Glick P, Xu J (2002) A model of (often mixed) stereotype content: Competence and warmth respectively follow from perceived status and competition. *J Pers Soc Psychol* 82:878–902.
- Henley NM (1989) Molehill or mountain? What we know and don't know about sex bias in language. *Gender and Thought: Psychological Perspectives*, eds Crawford M, Gentry M (Springer, New York), pp 59–78.
- Hellinger M, Bußmann H eds (2001) *Gender Across Languages: The Linguistic Representation of Women and Men*, IMPACT: Studies in Language and Society (John Benjamins Publishing Company, Amsterdam), Vol 9.
- Collobert R, et al. (2011) Natural language processing (almost) from scratch. *J Machine Learn Res* 12:2493–2537.
- Bolukbasi T, Chang KW, Zou JY, Saligrama V, Kalai AT (2016) Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems* 29, eds Lee DD, Sugiyama M, Luxburg UV, Guyon I, Garnett R (Curran Associates, Inc, Barcelona), pp 4349–4357.
- Caliskan A, Bryson JJ, Narayanan A (2017) Semantics derived automatically from language corpora contain human-like biases. *Science* 356:183–186.
- Zhao J, Wang T, Yatskar M, Ordóñez V, Chang KW (2017) Men also like shopping: Reducing gender bias amplification using corpus-level constraints. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, eds Palmer M, Hwa R (Association for Computational Linguistics, Copenhagen), pp 2979–2989.
- van Miltenburg E (2016) Stereotyping and bias in the Flickr30k dataset. arXiv: 1605.06083.
- Mikolov T, Chen K, Corrado G, Dean J (2013) Efficient estimation of word representations in vector space. arXiv:1301.3781.
- Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J (2013) Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems* 26, eds Burges CJ, Bottou L, Welling M, Ghahramani Z, Weinberger KQ (Curran Associates, Inc, Lake Tahoe, NV), pp 3111–3119.
- Hamilton WL, Leskovec J, Jurafsky D (2016) Diachronic word embeddings reveal statistical laws of semantic change. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, eds Erk K, Smith NA (Association for Computational Linguistics, Berlin), Vol 1, pp 1489–1501.
- Pennington J, Socher R, Manning C (2014) Glove: Global vectors for word representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, eds Moschitti A, Pang B (Association for Computational Linguistics, Doha, Qatar), pp 1532–1543.
- Sandhaus E (2008) *The New York Times Annotated Corpus* (Linguistic Data Consortium, Philadelphia).
- Ruggles S, Genadek K, Goeken R, Grover J, Sobek M (2015) Integrated Public Use Microdata Series: Version 6.0. Available at doi.org/10.18128/D010.V6.0. Accessed August 16, 2017.
- Osajima K (2005) Asian Americans as the model minority: An analysis of the popular press image in the 1960s and 1980s. *A Companion to Asian American Studies*, ed Ono KA (Blackwell Publishing Ltd, Malden, MA), pp 215–225.
- Fong TP (2002) *The Contemporary Asian American Experience: Beyond the Model Minority* (Prentice Hall, Upper Saddle River, NJ).
- Bryson V (2016) *Feminist Political Theory* (Palgrave Macmillan, New York).
- Rosen R (2013) *The World Split Open: How the Modern Women's Movement Changed America* (Tantor eBooks, Old Saybrook, CT).
- Thorne B, Henley N, Kramarae C, eds (1983) *Language, Gender, and Society* (Newbury House, Rowley, MA).
- Eckert P, McConnell-Ginet S (2003) *Language and Gender* (Cambridge Univ Press, Cambridge, UK).
- Evans S (2010) *Tidal Wave: How Women Changed America at Century's End* (Simon and Schuster, New York).
- Tasca C, Rapetti M, Carta MG, Fadda B (2012) Women and hysteria in the history of mental health. *Clin Pract Epidemiol Ment Health* 8:110–119.
- Sanghani R (2016) Feisty, frigid and frumpy: 25 Words we only use to describe women. The Telegraph. Available at https://www.telegraph.co.uk/women/life/ambitious-frigid-and-frumpy-25-words-we-only-use-to-describe-wom/. Accessed August 21, 2017.
- Zong J, Batalova J (2016) *Asian Immigrants in the United States* (Migration Policy Institute, Washington, DC).
- Lee SJ (1994) Behind the model-minority stereotype: Voices of high- and low-achieving Asian American students. *Anthropol Educ Q* 25:413–429.
- Kim A, Yeh CJ (2002) *Stereotypes of Asian American students* (ERIC Digest New York, NY).
- Lee SJ (2015) *Unraveling the "Model Minority" Stereotype: Listening to Asian American Youth* (Teachers College Press, New York), 2nd Ed.
- Gutiérrez DG (2016) A historic overview of Latino immigration and the demographic transformation of the United States. *The New Latino Studies Reader: A Twenty-First-Century Perspective*, eds Gutierrez RA, Almaguer T (Univ of California Press, Oakland, CA), pp 108–125.
- Levanon A, England P, Allison P (2009) Occupational feminization and pay: Assessing causal dynamics using 1950–2000 U.S. Census data. *Soc Forces* 88:865–891.
- Rothe S, Schütze H (2016) Word embedding calculus in meaningful ultradense subspaces. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, eds Erk K, Smith NA (Association for Computational Linguistics, Berlin), Vol 2, pp 512–517.
- Rudolph M, Ruiz F, Athey S, Blei D (2017) Structured Embedding Models for Grouped Data in Advances in Neural Information Processing Systems 30, eds Guyon I, et al. (Curran Associates, Inc, Long Beach, CA), pp 250–260.
- Rudolph M, Blei D (2017) Dynamic Bernoulli embeddings for language evolution. arXiv:1703.08052.
- Davies M (2010) The 400 million word corpus of historical American English (1810–2009). *Selected Papers from the Sixteenth International Conference on English Historical Linguistics (ICEHL 16)*, Pécs, 23–27 August 2010, eds Hegedüs I, Fodor A (John Benjamins Publishing, Amsterdam), Vol 325.
- Ullmann S (1962) *Semantics: An Introduction to the Science of Meaning* (Barnes & Noble, New York).
- Blank A (1999) *Why Do New Meanings Occur? A Cognitive Typology of the Motivations for Lexical Semantic Change in Historical Semantics and Cognition*, ed Koch P (Walter de Gruyter, New York).

51. Kulkarni V, Al-Rfou R, Perozzi B, Skiena S (2015) Statistically significant detection of linguistic change. *Proceedings of the 24th International Conference on World Wide Web*, eds Gangemi A, Leonardi S, Panconesi A (International World Wide Web Conferences Steering Committee, New York), pp 625–635.
52. Hamilton WL, Leskovec J, Jurafsky D (2016) Cultural shift or linguistic drift? comparing two computational measures of semantic change. *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, eds Su J, Duh K, Carreras X (Association for Computational Linguistics, Austin, TX), pp 2116–2121.
53. Chalabi M, Flowers A (2018) Dear Mona, what's the most common name in America? Available at <https://fivethirtyeight.com/features/whats-the-most-common-name-in-america/>. Accessed September 3, 2017.
54. Gunkel P (2013) 638 Primary personality traits. Available at ideonomy.mit.edu/essays/traits.html. Accessed August 21, 2017.