
A few of Dan Jurafsky's contributions to NLP

A brief introduction to the Stanford NLP group along with a few interesting papers co-authored by Dan Jurafsky

Dan Jurafsky

PhD in Computer Science

MacArthur Award

The Language of Food:
A Linguist Reads the
Menu



1992

2002

2014

2002

2012

2017

First automatic system
for semantic role labeling
with Daniel Gildea

Natural Language
Processing - Coursera

Predicting Sales from
the Language of Product
Descriptions

Semantic Role Labeling

Detection of the semantic arguments associated with the predicate or verb of a sentence and their classification into their specific roles

For instance , in : “*Dan sold the book to Daniel*”

The verb “to sell” is the *predicate*.

“Dan” is the seller or the *agent*

“The book” is a good being sold, or the *theme*

“Daniel” is the recipient.

The Language of Food

“ These cupcakes, they're like crack”

“Be warned, the wings are addicting”

“ Every time I need a fix, that fried chicken is so damn good.”

Somehow if it's a drug or we're addicted, it's really not really our fault. It's really the fault of the food which is this awful drug-like thing. It wasn't my fault. I had to eat that cupcake. It made me eat it.

(Shame on you cupcake)



Language from police body camera footage shows racial disparities in officer respect

Rob Voigt^{a,1}, Nicholas P. Camp^b, Vinodkumar Prabhakaran^c, William L. Hamilton^c, Rebecca C. Hetey^b, Camilla M. Griffiths^b, David Jurgens^c, Dan Jurafsky^{a,c}, and Jennifer L. Eberhardt^{b,1}

^aDepartment of Linguistics, Stanford University, Stanford, CA 94305; ^bDepartment of Psychology, Stanford University, Stanford, CA 94305; and ^cDepartment of Computer Science, Stanford University, Stanford, CA 94305

Contributed by Jennifer L. Eberhardt, March 26, 2017 (sent for review February 14, 2017; reviewed by James Pennebaker and Tom Tyler)

Significance

Police officers speak significantly less respectfully to black than to white community members in everyday traffic stops, even after controlling for officer race, infraction severity, stop location, and stop outcome. This paper presents a systematic analysis of officer body-worn camera footage, using computational linguistic techniques to automatically measure the respect level that officers display to community members. This work demonstrates that body camera footage can be used as a rich source of data rather than merely archival evidence, and paves the way for developing powerful language-based tools for studying and potentially improving police–community relations.

Study 1. Perceptions of Officer Treatment from Language

Study 2. Linguistic Correlates of Respect

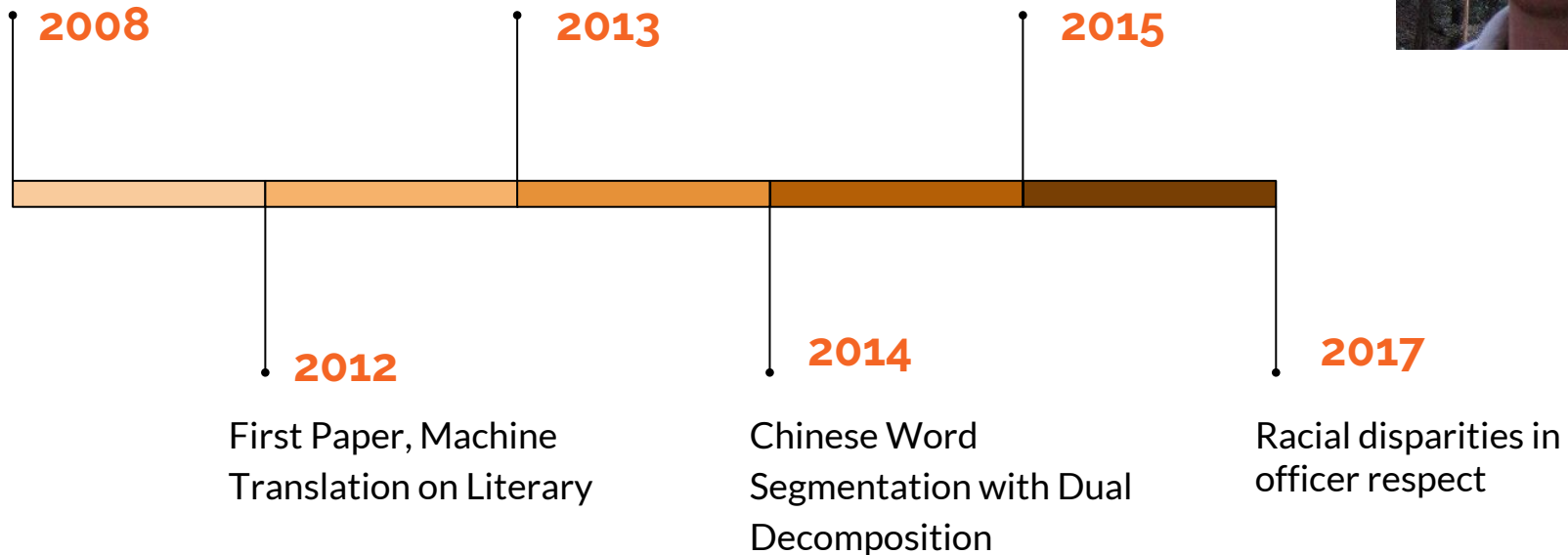
Study 3. Racial Disparities in Respect

Rob Voigt

BA, Chinese; Vassar
College

MA & PhD. Stanford Univ.
Chinese Poetry

The Users Who Say 'Ni':
Audience Identification
in Chinese-language
Restaurant Reviews



William L. Hamilton

Loyalty in Online
Communities. Reddit

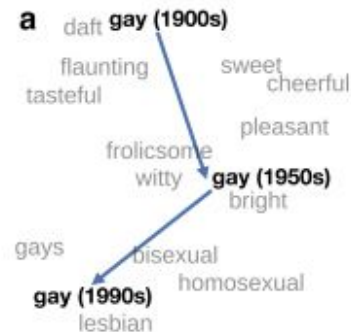
PhD. Stanford Univ. NLP

BA,MA; McGill Univ



[Modelling Sparse
Dynamical Systems with
Compressed Predictive
State Representations](#)

Diachronic Word
Embeddings Reveal
Statistical Laws of
Semantic Change



Study 1: Perceptions of Officer Treatment from Language.

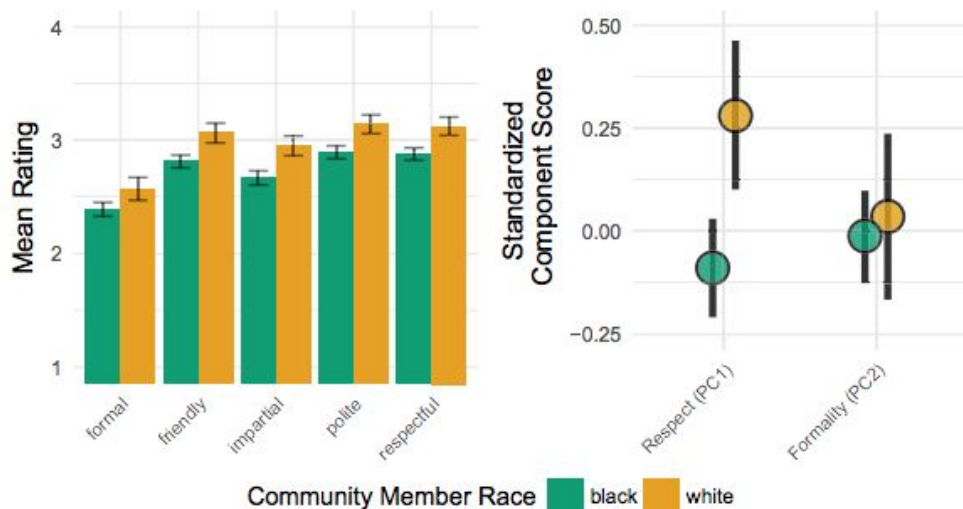


Fig. 1. (Left) Differences in raw participant ratings between interactions with black and white community members. (Right) When collapsed to two uncorrelated components, Respect and Formality, we find a significant difference for Respect but none for Formality. Error bars represent 95% confidence intervals. PC, principal component.

414 utterances; 312 Black and 102 White

	PC1: RESPECT	PC2: FORMALITY
Formal	0.272	0.913
Friendly	0.464	-0.388
Impartial	0.502	-0.113
Polite	0.487	-0.047
Respectful	0.471	0.026
% of Variance Explained	71.3%	21.9%

Drawbacks:
 Scale, 26 million stops per year
 Sample size, 414, too small

Study 2: Linguistic Correlates of Respect.

Feature Name	Implementation	Source
Adverbial "Just"	"Just" occurs in a dependency arc as the head of an <i>advmod</i> relation	
Apologizing	Lexicon: "sorry", "oops", "woops", "excuse me", "forgive me", "apologies", "apologize", "my bad", "my fault"	[4]
Ask for Agency	Lexicon: "do me a favor", "let me", "allow me", "can i", "should i", "may i", "might i", "could i"	[4]
Bald Command	The first word in a sentence is a bare verb with part-of-speech tag VB ("look", "give", "wait" etc.) but is not one of "be", "do", "have", "thank", "please", "hang".	
Colloquialism	Regular expression capturing "y'all", "ain't" and words ending in "in'" such as "walkin'", "talkin'", etc., as marked by transcribers	
Conditional	Lexicon: "if"	
Disfluency	Word fragment ("Well I thi-") as indicated by transcribers	[5, 6]
Filled Pauses	Lexicon: "um", "uh"	[7, 8]
First Names	Top 1000 most common first names from the 1990 US Census, where first letter is capitalized in transcript	[9, 10] ¹
Formal Titles	Lexicon: "sir", "ma'am", "maam", "mister", "mr*", "ms*", "madam", "miss", "gentleman", "lady"	[9, 10]
For Me	Lexicon: "for me"	
For You	Lexicon: "for you"	
Give Agency	Lexicon: "let you", "allow you", "you can", "you may", "you could"	[4]
Gratitude	Lexicon: "thank", "thanks", "appreciate"	[4]
Goodbye	Lexicon: "goodbye", "bye", "see you later"	
Hands on the Wheel	Regular expression capturing cases like "keep your hands on the wheel" and "leave your hands where I can see them": "hands? (<[.,?!:;]+)?(wheel see)"	
Hedges	All words in the "Tentat" LIWC lexicon	[11]
Impersonal Pronoun	All words in the "Imppron" LIWC lexicon	[4, 11]
Informal Titles	Lexicon: "dude*", "bro*", "boss", "bud", "buddy", "champ", "man", "guy", "guy", "brotha", "sista", "son", "sonny", "chief"	[9, 10, 12]
Introductions	Regular expression capturing cases like "I'm Officer [name] from the OPD" and "How's it going?": "((i my name).+officer officer.+(oak and opd))! (hi helle hey good afternoon good morning good evening how are you doing how 's it going))"	[4]
Last Names	Top 5000 most common last names from the 1990 US Census, where first letter is capitalized in transcript	[9, 10] ²
Linguistic Negation	All words in the "Negate" LIWC lexicon	[11]
Negative Words	All words in the "Negativ" category in the Harvard General Inquirer, matching on word lemmas	[4, 13]
Positive Words	All words in the "Positiv" category in the Harvard General Inquirer, matching on word lemmas	[4, 13]

EXAMPLE	RESPECT SCORE
<p>FIRST NAME ASK FOR AGENCY QUESTIONS</p> <p>[name], can I see that driver's license again?</p> <p>It- it's showing suspended. Is that- that's you?</p> <p>DISFLUENCY NEGATIVE WORD DISFLUENCY</p>	-1.07
<p>INFORMAL TITLE ASK FOR AGENCY ADVERBIAL "JUST"</p> <p>All right, my man. Do me a favor. Just keep your hands on the steering wheel real quick.</p> <p>"HANDS ON THE WHEEL"</p>	-0.51
<p>APOLOGY INTRODUCTION LAST NAME</p> <p>Sorry to stop you. My name's Officer [name] with the Police Department.</p>	0.84
<p>FORMAL TITLE SAFETY PLEASE</p> <p>There you go, ma'am. Drive safe, please.</p>	1.21
<p>ADVERBIAL "JUST" FILLED PAUSE REASSURANCE</p> <p>It just says that, uh, you've fixed it. No problem.</p> <p>Thank you very much, sir.</p> <p>GRATITUDE FORMAL TITLE</p>	2.07

Study 2: Linguistic Correlates of Respect.

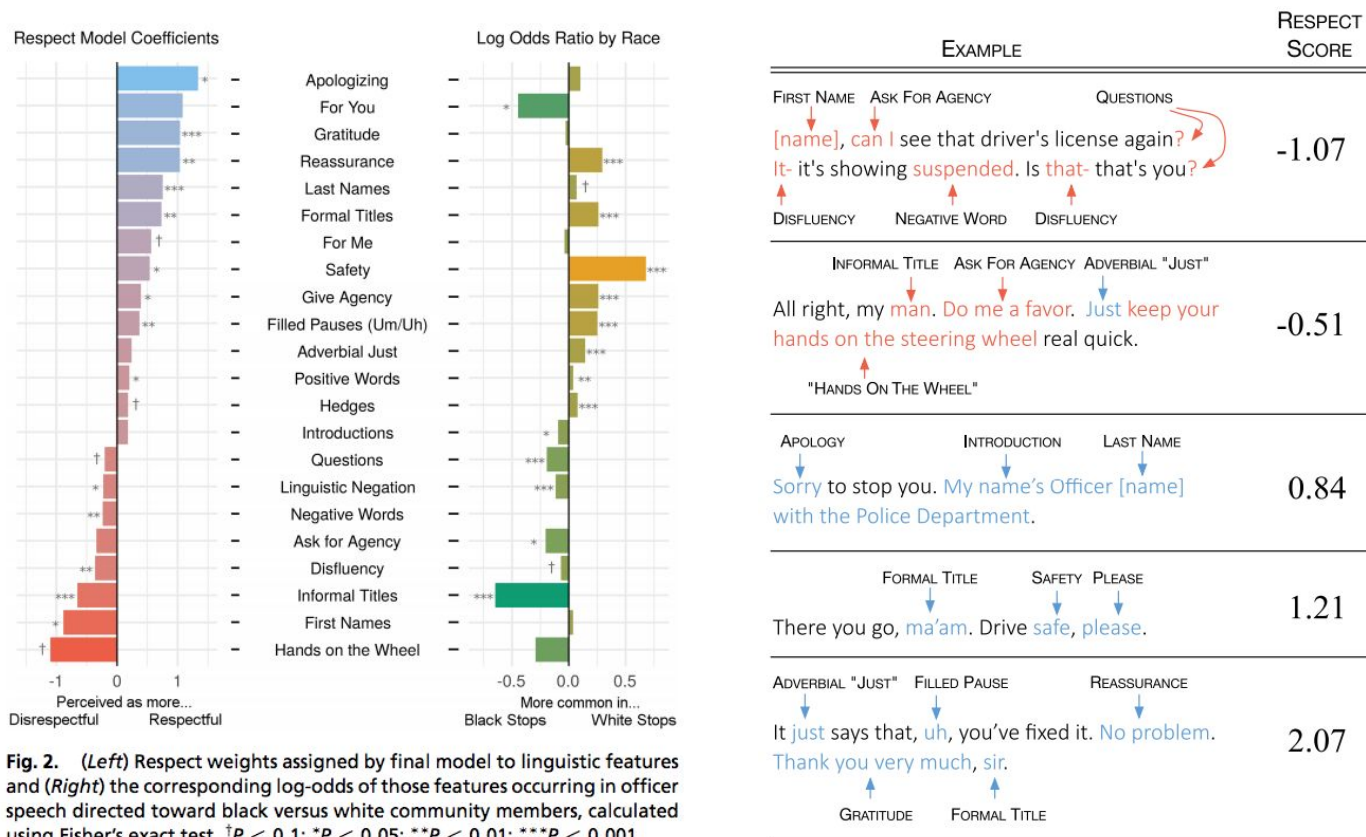


Fig. 2. (Left) Respect weights assigned by final model to linguistic features and (Right) the corresponding log-odds of those features occurring in officer speech directed toward black versus white community members, calculated using Fisher's exact test. † $P < 0.1$; * $P < 0.05$; ** $P < 0.01$; *** $P < 0.001$.

Study 3: Racial Disparities in Respect.

- 36,738 utterances
- community member race, age, and gender
- officer race
- whether a search was conducted
- the result of the stop (warning, citation, or arrest)

	<i>Respect</i>			<i>Formality</i>		
	β	CI	p	β	CI	p
Fixed Parts						
Arrest Occurred	0.00	-0.03 – 0.03	.933	0.01	-0.02 – 0.04	.528
Citation Issued	0.04	0.02 – 0.06	<.001	0.01	-0.01 – 0.03	.209
Search Conducted	-0.08	-0.11 – -0.05	<.001	0.00	-0.03 – 0.02	.848
Age	0.07	0.05 – 0.09	<.001	0.05	0.03 – 0.07	<.001
Gender (F)	0.02	0.00 – 0.04	.062	0.02	0.00 – 0.04	.025
Race (W)	0.05	0.03 – 0.08	<.001	-0.01	-0.04 – 0.01	.236
Officer Race (B)	0.00	-0.03 – 0.04	.884	0.00	-0.03 – 0.03	.987
Officer Race (O)	0.00	-0.04 – 0.03	.809	0.00	-0.03 – 0.02	.783
Officer Race (B) : Race (W)	-0.01	-0.03 – 0.02	.583	0.01	-0.01 – 0.03	.188
Officer Race (O) : Race (W)	-0.01	-0.03 – 0.02	.486	0.00	-0.02 – 0.02	.928

Study 3: Racial Disparities in Respect.

- Other hypothesis
- Are the racial disparities in the respectfulness of officer speech we observe driven by a small number of officers? **NO!**

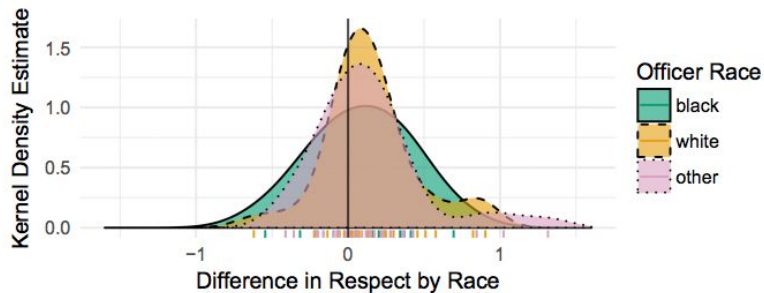


Fig. 4. Kernel density estimate of individual officer-level differences in Respect when talking to white as opposed to black community members, for the 90 officers in our dataset who have interactions with both blacks and whites. More positive numbers on the x axis represent a greater positive shift in Respect toward white community members.

Study 3: Racial Disparities in Respect.

- Prediction

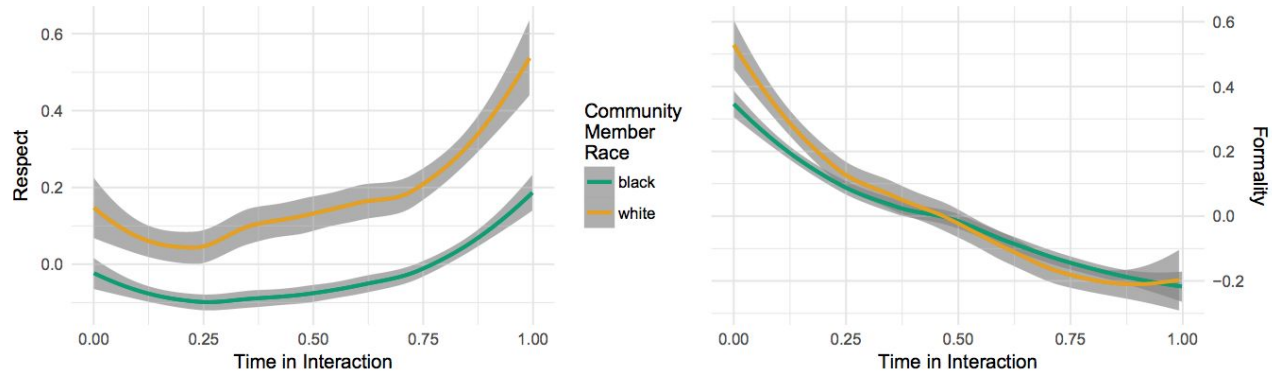


Fig. 5. Loess-smoothed estimates of the (Left) Respect and (Right) Formality of officers' utterances relative to the point in an interaction at which they occur. Respect tends to start low and increase over an interaction, whereas the opposite is true for Formality. The race discrepancy in Respect is consistent throughout the interactions in our dataset.

Summary

First time researchers use body-worn camera footage to explore racial disparities in officer's respect towards Black and White community members.

Significant racial disparities are found, but the causes of the disparities is not clear

Reid Pryzant

BA, Computer Science and
Biology at Williams College



2016

2016 - Present

PhD in Computer
Science at Stanford
University

2017

Predicting Sales from
the Language of Product
Descriptions

Predicting Sales from the Language of Product Descriptions

Predicting Sales from the Language of Product Descriptions

Reid Pryzant
Stanford University
rpryzant@stanford.edu

Young-joo Chung
Rakuten Institute of Technology
yjchung@acm.org

Dan Jurafsky
Stanford University
jurafsky@stanford.edu

ABSTRACT

What can a business say to attract customers? E-commerce vendors frequently sell the same items but use different marketing strategies to present their goods. Understanding consumer responses to this heterogeneous landscape of information is important both as business intelligence and, more broadly, a window into consumer attitudes. When studying consumer behavior, the existing literature is primarily concerned with product reviews. In this paper we posit that textual product descriptions are also important determinants of consumer choice. We mine 90,000+ product descriptions on the Japanese e-commerce marketplace Rakuten and identify actionable writing styles and word usages that are highly predictive of consumer purchasing behavior. In the process, we observe the inadequacies of traditional feature extraction algorithms, namely their inability to control for the implicit effects of confounds like brand loyalty and pricing strategies. To circumvent this problem, we propose a novel neural network architecture that leverages an adversarial objective to control for confounding factors, and attentional scores over its input to automatically elicit textual features as a domain-specific lexicon. We show that these textual features can predict the sales of each product, and investigate the narratives highlighted by these words. [Our results suggest that appeals to authority, polite language, and mentions of informative and seasonal language win over the most customers.](#)

test, and evaluate products before making purchasing decisions, the remote nature of e-commerce renders such tactile evaluations obsolete.

In lieu of in-store evaluation, online shoppers increasingly rely on alternative sources of information. This includes “word-of-mouth” recommendations from outside sources [9] and local product reviews [13, 18, 20]. These factors, though well studied, are only indirectly controllable from a business perspective [25, 52]. Business owners have considerably stronger control over their own product descriptions. The same products may be sold by multiple vendors, with each item having a different textual description (note that we take *product* to mean a purchasable object, and *item* to mean an individual e-commerce listing). Studying consumers’ reactions to these descriptions is valuable both as business intelligence and as a new window into consumer attitudes.

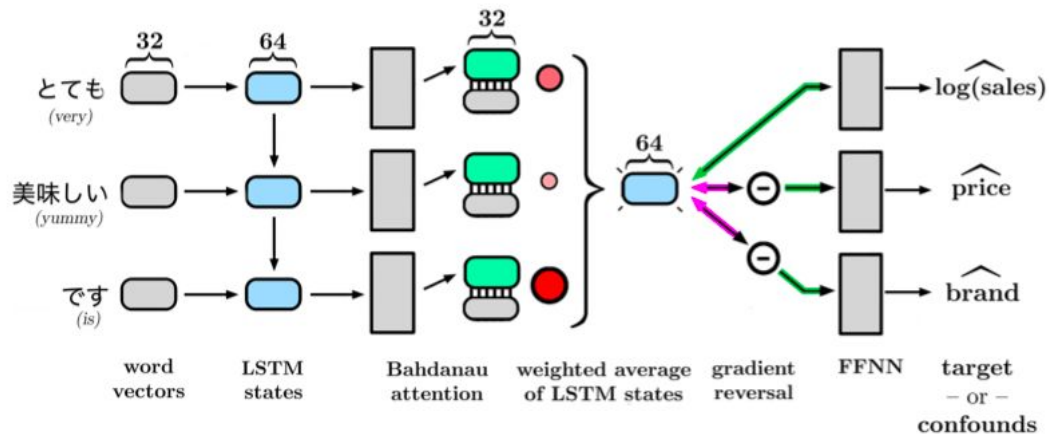
The hypothesis that business-generated product descriptions affect consumer behavior (manifested in sales) has received strong support in prior empirical studies [22, 26, 34, 37, 39]. However, these studies have only used summary statistics of these descriptions (i.e. readability, length, completeness). We propose that embedded in these product descriptions are narratives that affect shoppers, which can be studied by examining the words in each description.

Our hypothesis is that product descriptions are fundamentally a kind of social discourse, one whose linguistic contents have real control over consumer purchasing behavior. Business owners em-

Observations leading up to the paper

- Human judgment and behavior is influenced by persuasive rhetoric
- Business owners employ narratives to portray their products, and consumers react accordingly according to their beliefs and attitudes
- Aim to unearth actionable phrases that can help e-commerce vendors increase their sales regardless of what's being sold
- We wish to study the impact of linguistic structures in product descriptions in isolation, beyond those indicators of price or branding.

Proposed Model



- Forward Pass where predictions are generated
- Backward Pass where parameters are updated
- Feature Selection using attentional scores

Influential words

1. Informativeness
2. Authority
3. Seasonality
4. Politeness

Two product descriptions of the same product

Royce's chocolate has become a **standard** Hokkaido **souvenir**. They are packaged one by one so your hands won't get dirty! Also, our **staff** recommends this product!

vs

Four types of nuts: almonds, cashews, pecans, macadamia, as well as cookie crunch and almond puff were packed carefully into each chocolate bar. This item is shipped with a refrigerated courier service during the **summer**.

Summary

- Hypothesis is that product descriptions are fundamentally a kind of social discourse, one whose linguistic contents have real control over consumer purchasing behavior
- Used Deep Adversarial Feature Mining based model
- Influential words

Interesting Reads

1. The Language Of Food : <https://web.stanford.edu/~jurafsky/thelanguageoffood.html>
2. Dan Jurafsky's Blog - <http://languageoffood.blogspot.com>
3. Loyalty in Online Communities: https://www.cs.cornell.edu/~cristian/index_files/loyalty.pdf

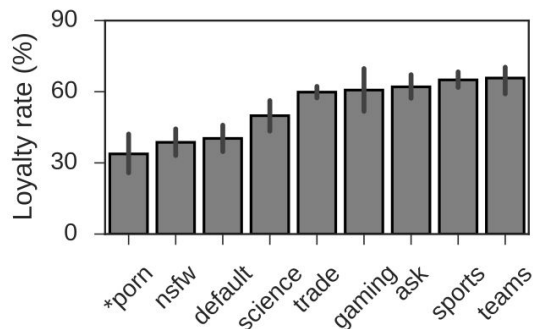


Figure 1: **Loyalty rates by community category.** Most categories were scraped from /r/ListOfSubreddits, while the category labels containing “*” were generated by matching on subreddit names fitting the specified pattern. Note that “*porn” are image-sharing communities like /r/EarthPorn, not pornography; the “nsfw” category contains pornography and other explicit content. 99% bootstrapped CIs are shown.

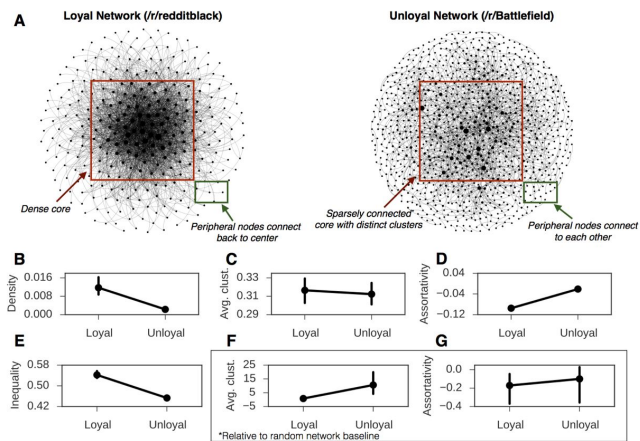


Figure 2: **Loyalty reflected in community interaction networks.** Networks of user interactions reveal how loyalty manifests at the collective, community level. **A**, Example loyal and unloyal networks built from interactions in March, 2014 in two war-related video game communities. **B-D**, Differences in network statistics for the empirical networks. **E**, Difference in activity inequality (measured using the Gini coefficient). **F-G**, Network statistics relative to a null configuration model baseline. Bootstrapped 99.9% confidence intervals are shown.

Thank You!