

When do Words Matter? Understanding the Impact of Lexical Choice on Audience Perception using Individual Treatment Effect Estimation

Zhao Wang and Aron Culotta

Department of Computer Science
Illinois Institute of Technology, Chicago, IL 60616
zwang185@hawk.iit.edu, aculotta@iit.edu

Abstract

Studies across many disciplines have shown that lexical choice can affect audience perception. For example, how users describe themselves in a social media profile can affect their perceived socio-economic status. However, we lack general methods for estimating the causal effect of lexical choice on the perception of a specific sentence. While randomized controlled trials may provide good estimates, they do not scale to the potentially millions of comparisons necessary to consider all lexical choices. Instead, in this paper, we first offer two classes of methods to estimate the effect on perception of changing one word to another in a given sentence. The first class of algorithms builds upon quasi-experimental designs to estimate individual treatment effects from observational data. The second class treats treatment effect estimation as a classification problem. We conduct experiments with three data sources (Yelp, Twitter, and Airbnb), finding that the algorithmic estimates align well with those produced by randomized-control trials. Additionally, we find that it is possible to transfer treatment effect classifiers across domains and still maintain high accuracy.

1 Introduction¹

Numerous examples from cognitive science, linguistics, and marketing show that lexical choice can affect audience perception (Danescu-Niculescu-Mizil et al. 2012; Ludwig et al. 2013; Thibodeau and Boroditsky 2013; Riley and Luippold 2015; Reddy and Knight 2016; Preoŝiuc-Pietro, Guntuku, and Ungar 2017; Packard and Berger 2017; Nguyen et al. 2017). For example, a social media user who writes “*I’m excited!*” may be more likely to be perceived as female than one who writes “*I’m stoked!*” (Reddy and Knight 2016). Similarly, a book with a review containing the sentence “*I loved this book!*” may be perceived as more desirable than one with a review stating “*An excellent novel.*” (Ludwig et al. 2013).

Despite this prior work, we still lack general methods for estimating the causal effect on perception of a single linguistic change in a specific sentence. For example, how much

does changing the word “*excited*” to “*stoked*” in the example above increase the chance that a reader will infer the user to be male? Being able to answer such questions has implications not only for marketing and public messaging campaigns, but also for author obfuscation (Hagen, Pothast, and Stein 2017) and stylistic deception detection (Afroz, Brennan, and Greenstadt 2012).

A standard empirical approach is to conduct a Randomized Control Trial (RCT), in which subjects are shown texts that differ only in a single linguistic change, and are subsequently asked to rate their perception with respect to a particular attribute. By controlling for the context, we can then attribute changes in perception to the single linguistic change.

Unfortunately, it is impractical to scale such RCTs to the many possible word substitutions across thousands of sentences, making applications based on such methods infeasible. The goal of this paper is to instead investigate automated methods that estimate how a specific lexical choice affects perception of a single sentence. Our approach builds upon a type of causal inference called *Individual Treatment Effect* (ITE) estimation. An ITE estimation algorithm estimates the effect of an intervention on an individual; e.g., how effective a drug will be for a specific person. Recently, a number of ITE estimators have been proposed that require only observational data, based on Rubin’s potential outcome framework (Rubin 1974). In this paper, we formulate our problem as a type of ITE estimation, which we call *Lexical Substitution Effect* (LSE) estimation. We propose two classes of LSE estimators. The first class adapts previous algorithms in ITE estimation to the task of LSE estimation. These methods take as input sentences labeled according to attributes of interest (e.g., a tweet labeled by the gender of the author) and then produces tuples of the form $(w_i, w_j, s, \hat{\tau})$, indicating the estimated LSE ($\hat{\tau}$) of changing the word w_i to w_j for sentence s , with respect to the attribute of interest. The second class of estimator is inspired by recent work that frames causal inference as a classification problem (Lopez-Paz et al. 2015). This approach requires some labeled examples of the form (w_i, w_j, s, τ) , where τ is the “true” LSE according to a RCT. It then fits a classifier based on properties of (w_i, w_j, s) to produce LSE estimates for new sentences.

We conduct studies using three data sources: Airbnb listings, Yelp reviews, and Twitter messages. For Airbnb, we

Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹An expanded version of this paper is available at: <https://arxiv.org/abs/1811.04890> ; replication files and data are available at: <https://github.com/tapilab/aaai-2019-words>

consider the perception of the desirability of a rental based on a sentence from the listing description. For Yelp and Twitter, we consider the perception of the gender of the author. We estimate LSE for thousands of word substitutions across millions of sentences, comparing the results of different LSE estimators. For a sample of sentences, we additionally conduct RCTs using Amazon Mechanical Turk to validate the quality of the algorithmic estimates with respect to human judgments. Overall, we find that the algorithmic estimates align well with those produced by RCTs. We also find that it is possible to transfer treatment effect classifiers across domains and still maintain high quality estimates.

2 Related Work

Studies investigating the effect of wording in communication strategies dates back at least 60 years (Hovland, Janis, and Kelley 1953). Recent research has explored the effect of wording on Twitter message propagation (Tan, Lee, and Pang 2014), how word choice and sentence structure affects memorability of movie quotes (Danescu-Niculescu-Mizil et al. 2012), and how characteristics of news articles influence with high story sharing rates (Berger and Milkman 2012). Additionally, there has been recent psycho-linguistic research discovering how to infer user attributes (e.g., gender, age, occupation) based on language styles. Preotiuc-Pietro, Xu, and Ungar (2016) explore a wide set of meaningful stylistic paraphrase pairs and verified a number of psycho-linguistic hypotheses about the effect of stylistic phrase choice on human perception. Preotiuc-Pietro, Guntuku, and Ungar (2017) further conduct experiment to control human perception of user trait in tweets. Similarly, Reddy and Knight (2016) propose methods to obfuscate gender by lexical substitutions.

As a type of causal inference, individual treatment effect estimation is typically explored in medical trials to estimate effects of drug use on a health outcome. Classical approaches include nearest-neighbor matching, kernel methods and so on (Crump et al. 2008; Lee 2008; Willke et al. 2012). However, the performance of these methods do not scale well with the number of covariates (Wager and Athey 2017). To accommodate a large number of complex covariates, researchers have recently explored techniques such as random forests (Breiman 2001) and causal forests (Wager and Athey 2017). Motivated by their successful applications in the medical domain, we propose to adapt these techniques to the linguistic domain. Specifically, we conduct experiments to algorithmically estimate the causal effect of lexical change on perception for a single sentence.

In summary, while some prior work has studied overall effects of lexical substitution, in this paper we instead propose methods to estimate context-specific effects. That is, we are interested in quantifying the effect on perception caused by a single word change in a specific sentence. The primary contributions are (1) to formalize the LSE problem as a type of ITE; (2) to adapt ITE methods to the text domain; (3) to develop classifier-based estimators that are able to generalize across domains.

3 Individual Treatment Effect Estimation

In this section, we first provide background on Individual Treatment Effect (ITE) estimation, and then in the following section we will adapt ITE to Lexical Substitution Effect (LSE) estimation.

Assume we have dataset D consisting of n observations $D = \{(\mathbf{X}_1, T_1, Y_1), \dots, (\mathbf{X}_n, T_n, Y_n)\}$, where \mathbf{X}_i is the *covariate vector* for individual i , $T_i \in \{0, 1\}$ is a binary *treatment indicator* representing whether i is in the treatment ($T_i = 1$) or control ($T_i = 0$) group, and Y_i is the observed *outcome* for individual i . For example, in a pharmaceutical setting, i is a patient; \mathbf{X}_i is a vector of the socio-economic variables (e.g., gender, age, height); T_i indicates whether he did ($T_i = 1$) or did not ($T_i = 0$) receive the medication treatment, and $Y_i \in \{0, 1\}$ indicates whether he is healthy ($Y_i = 1$) or sick ($Y_i = 0$).

We are interested in quantifying the causal effect that treatment T has on the outcome Y . The fundamental problem of causal inference is that we can only observe one outcome per individual, either the outcome of an individual receiving a treatment or not. Thus, we do not have direct evidence of what might have happened had we given individual i a different treatment. Rubin’s potential outcome framework is a common way to formalize this fundamental problem (Rubin 1974). Let $Y^{(1)}$ indicate the potential outcome an individual would have got had they received treatment ($T = 1$), and similarly let $Y^{(0)}$ indicate the outcome an individual would have got had they received no treatment ($T = 0$). While we cannot observe both $Y^{(1)}$ and $Y^{(0)}$ at the same time, we can now at least formally express several quantities of interest. For example, we are often interested in the *average treatment effect* (τ), which is the expected difference in outcome had one received treatment versus not: $\tau = \mathbb{E}[Y^{(1)}] - \mathbb{E}[Y^{(0)}]$. In this paper, we are interested in the *Individual Treatment Effect* (ITE), which is the expected difference in outcome for a specific type of individual:

$$\tau(\mathbf{x}) = \mathbb{E}[Y^{(1)} | \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y^{(0)} | \mathbf{X} = \mathbf{x}] \quad (1)$$

that is, the treatment effect for individuals where $\mathbf{X} = \mathbf{x}$. For example, if the covariate vector represents the (age, gender, height) of a person, then the ITE will estimate treatment effects for individuals that match along those variables.

Estimating $\tau(\mathbf{x})$ from observational data, in which we have no control over the treatment assignment mechanism, is generally intractable due to the many possible confounds that can exist (e.g., patients receiving the drug may be *a priori* healthier on average than those not receiving the drug). However, numerous algorithms exist to produce estimates of $\tau(\mathbf{x})$ from observational data, for example propensity score matching (Austin 2008). These methods require additional assumptions, primarily the *Strongly Ignorable Treatment Assignment* (SITA) assumption. SITA assumes that the treatment assignment is conditionally independent of the outcome given the covariate variables: $T \perp \{Y^{(0)}, Y^{(1)}\} | \mathbf{X}$. While this assumption does not hold generally, methods built on this assumption have often been found to work well.

With SITA, we can estimate ITE using only observational

data as follows:

$$\hat{\tau}(\mathbf{x}) = \mathbb{E}[Y|T = 1, \mathbf{X} = \mathbf{x}] - \mathbb{E}[Y|T = 0, \mathbf{X} = \mathbf{x}] \quad (2)$$

$$= \frac{1}{|S_1(\mathbf{x})|} \sum_{i \in S_1(\mathbf{x})} Y_i - \frac{1}{|S_0(\mathbf{x})|} \sum_{i \in S_0(\mathbf{x})} Y_i \quad (3)$$

where $S_1(\mathbf{x})$ is the set of individuals i such that $\mathbf{X}_i = \mathbf{x}$ and $T_i = 1$, and similarly for $S_0(\mathbf{x})$. In other words, Equation (3) simply computes, for all individuals with covariates equal to \mathbf{x} , the difference between the average outcome for individuals in the treatment group and the average outcome for individuals in the control group. For example, if $\mathbf{X} = \mathbf{x}$ indicates individuals with (age=10, gender=male, height=5), $T = 1$ indicates that an individual receives drug treatment and $T = 0$ that they do not, then $\hat{\tau}(\mathbf{x})$ is the difference in average outcome between individuals who receive treatment and those who do not.

A key challenge to using Equation (3) in practice is that \mathbf{X} may be high dimensional, leading to a small sample where $\mathbf{X} = \mathbf{x}$. In the extreme case, there may be exactly one instance where $\mathbf{X} = \mathbf{x}$. Below, we describe several approaches to address this problem, which we will subsequently apply to LSE estimation tasks.

4 Lexical Substitution Effect Estimation

In this section, we apply concepts from §3 to estimate lexical substitution effect on perception. As a motivating example, consider the following two hypothetical sentences describing the neighborhood of an apartment listed on Airbnb:

A: There are plenty of **shops** nearby.

B: There are plenty of **boutiques** nearby.

We are interested in how substituting *shops* with *boutiques* affects the perceived desirability of the rental. E.g., because *boutiques* connotes a high-end shop, the reader may perceive the rental to be in a better neighborhood, and thus more desirable. Critically, we are interested in the effect of this substitution *in one particular sentence*. For example, consider a third sentence:

C: You can take a 10 minute ride to visit some **shops**.

We would expect the effect of substituting *shops* to *boutiques* in sentence **C** to be less than the effect for sentence **A**, since the word *shops* in **C** is less immediately associated with the rental.

First of all, to map the notation of §3 to this problem, we specify a sentence to be our primary unit of analysis (i.e., the “individual”). We make this choice in part for scalability and in part because of our prior expectation on effect sizes — it seems unlikely that a single word change will have much effect on the perception of a 1,000 word document, but it may affect the perception of a single sentence. The covariate vector \mathbf{X} represents the other words in a sentence, excluding the one being substituted. E.g., in example sentence **A**, $\mathbf{X} = \langle \text{There, are, plenty, of, } \rightarrow, \text{ nearby} \rangle$.

Second, we note that there are many possible lexical substitutions to consider for each sentence. If we let p index a substitutable word pair (*control word* \rightarrow *treatment word*), then we can specify T_i^p to be the lexical substitution assignment variable for sentence i . For example, if p represents

the substitutable word pair (*shops, boutiques*), then $T_i^p = 0$ indicates that sentence i is in the control group that has the control word *shops* in it, and we call it the control sentence, and $T_i^p = 1$ indicates that sentence i is treated by substituting the control word *shops* to the treatment word *boutiques*.

Third, the outcome variable Y indicates the perception with respect to a particular aspect (i.e., desirability or gender in this paper). For example, in the experiments below, we let $Y \in \{1, 2, 3, 4, 5\}$ be an ordinal variable expressing the perceived desirability level of an apartment rental based on a single sentence.

Finally, with these notations, we can then express the *Lexical Substitution Effect* (LSE), which can be understood as the ITE of performing the word substitution indicated by word pair p on a sentence with context words $\mathbf{X} = \mathbf{x}$:

$$\tau(\mathbf{x}, p) = \mathbb{E}[Y^{p(1)}|\mathbf{X} = \mathbf{x}] - \mathbb{E}[Y^{p(0)}|\mathbf{X} = \mathbf{x}] \quad (4)$$

If we have data of the form $D = \{(\mathbf{X}_1, T_1, Y_1), \dots, (\mathbf{X}_n, T_n, Y_n)\}$, we can then use the SITA assumption to calculate the LSE:

$$\hat{\tau}(\mathbf{x}, p) = \frac{1}{|S_1^p(\mathbf{x})|} \sum_{i \in S_1^p(\mathbf{x})} Y_i - \frac{1}{|S_0^p(\mathbf{x})|} \sum_{i \in S_0^p(\mathbf{x})} Y_i \quad (5)$$

where $S_1^p(\mathbf{x})$ is the set of sentences i such that $\mathbf{X}_i = \mathbf{x}$ and $T_i^p = 1$, and similarly for $S_0^p(\mathbf{x})$.

As mentioned previously, the high dimensionality of \mathbf{X} is the key problem with using Equation (5). This problem is even more critical in the linguistic domain than in traditional ITE studies in clinical domains — the total number of unique words is likely to be greater than the space of all possible socio-economic variables of a patient. For example, it is entirely possible that exactly one sentence in a dataset has context $\mathbf{X}_i = \mathbf{x}$.

In the subsections that follow, we first describe four algorithms from ITE estimation literature and how we adapt them to LSE estimation. Then, we describe a classifier-based approach that uses a small amount of labeled data to produce LSE estimates. As a running example, we will consider changing *shops* to *boutiques* in the sentence “*There are plenty of shops nearby.*” Sentences that contain the control word (e.g., “*shops*”) are called *control samples*, and those containing the treatment word (e.g., “*boutiques*”) are called *treatment samples*. Finally, since we only estimate LSE for one word substitution in one particular sentence each time, we will drop notation p in the following formulas.

4.1 K-Nearest Neighbor (KNN) Matching

KNN is a classical approach for non-parametric treatment effect estimation using nearest neighbor voting. The dimensionality problem is addressed by averaging the outcome variables of K closest neighbors. ITE estimation with KNN computes the difference between the average outcome of K nearest neighbors in treatment samples and control samples:

$$\hat{\tau}_{KNN}(\mathbf{x}) = \left(\frac{1}{K} \sum_{i \in S_1(\mathbf{x}, K)} Y_i \right) - \left(\frac{1}{K} \sum_{i \in S_0(\mathbf{x}, K)} Y_i \right) \quad (6)$$

For an individual with covariate $X = x$, $S_1(\mathbf{x}, K)$ and $S_0(\mathbf{x}, K)$ are the K nearest neighbors in treatment ($T = 1$) and control ($T = 0$) samples, respectively.

In LSE estimation with KNN matching, we first represent each sentence using standard tf-idf bag-of-words features, then apply cosine similarity to identify the K closest neighbor-sentences. For our running example, we get $S_0(\mathbf{x}, K)$ by selecting the K sentences that have highest cosine similarity with “*There are plenty of ... nearby*” from the control samples, and get set $S_1(\mathbf{x}, K)$ by selecting the K closest sentences to the treatment samples. Then, the KNN estimator calculates LSE by computing the difference between the average label values of K nearest sentences in the treatment samples and control samples.

4.2 Virtual Twins Random-Forest (VT-RF)

The virtual twins approach (Foster, Taylor, and Ruberg 2011) is a two step procedure. First, it fits a random forest with all observational data (including control samples and treatment samples), where each data is represented by inputs (\mathbf{X}_i, T_i) and outcome Y_i . Then, to estimate the ITE, it computes the difference between the predicted values for treatment input $(\mathbf{X}_i, T_i = 1)$ and control input $(\mathbf{X}_i, T_i = 0)$. The name ‘*virtual twin*’ derives from the fact that for the control input $(\mathbf{X}_i, T_i = 0)$, we make a copy $(\mathbf{X}_i, T_i = 1)$ as treatment input that is alike in every way to the control input except for the treatment variable. If $\hat{Y}(\mathbf{x}, 1)$ is the value predicted by the random forest for input $(\mathbf{X} = \mathbf{x}, T = 1)$, then the virtual twin estimate is:

$$\hat{\tau}_{VT}(\mathbf{x}) = \hat{Y}(\mathbf{x}, 1) - \hat{Y}(\mathbf{x}, 0) \quad (7)$$

where $\hat{Y}(\mathbf{x}, 1)$ is the outcome for the ‘*virtual twin*’ (treatment) input and $\hat{Y}(\mathbf{x}, 0)$ is the outcome for control input.

In LSE estimation with VT-RF, we first represent each sentence using binary bag-of-words features (which we found to be more effective than tf-idf). We then fit a random forest to estimate LSE by taking the difference in the posterior probabilities for the virtual twin sentence and the original sentence. For our running example, we fit a random forest classifier using all sentences containing either *shops* or *boutiques* except the current sentence (for out-of-bag estimation). Meanwhile, we generate the virtual twin sentence “*There are plenty of boutiques nearby.*” Then the estimated LSE is computed by taking the difference between $P(Y = 1 | \text{“There are plenty of boutiques nearby”})$ and $P(Y = 1 | \text{“There are plenty of shops nearby”})$.

4.3 Counterfactual Random-Forest (CF-RF)

Counterfactual random forest (Lu et al. 2018) is similar to VT-RF in that they both calculate ITE by taking the difference between predictions of random forest models. However, CF-RF is different from VT-RF by fitting two separate random forests: a control forest fitted with control samples, and a treatment forest fitted with treatment samples. The ITE is then estimated by taking the difference between the prediction (by treatment forest) for a treatment input

($\hat{Y}_1(\mathbf{x}, 1)$) and the prediction (by control forest) for a control input ($\hat{Y}_0(\mathbf{x}, 0)$):

$$\hat{\tau}_{CF}(\mathbf{x}) = \hat{Y}_1(\mathbf{x}, 1) - \hat{Y}_0(\mathbf{x}, 0) \quad (8)$$

In LSE estimation with CF-RF, after representing each sentence with binary bag-of-words features, we first fit a control random forest and a treatment random forest and then estimate LSE by taking probability difference between virtual twin sentence and the control sentence. For example, we fit a control forest with all sentences containing *shops* excluding the current one (for out-of-bag estimation) and a treatment forest with all sentences containing *boutiques*. We then estimate LSE by taking the difference between $P(Y = 1 | \text{“There are plenty of boutiques nearby”})$ predicted by treatment forest and $P(Y = 1 | \text{“There are plenty of shops nearby”})$ predicted by control forest.

4.4 Causal Forest (CSF)

A causal forest (Wager and Athey 2017) is a recently introduced model for causal estimation. While it also uses random forests, it modifies the node splitting rule to consider treatment heterogeneity. Whereas random forests create splits to maximize the purity of Y labels, causal forests instead create splits by maximizing the variance of estimated treatment effects in each leaf. To estimate ITE for an instance i , a causal forest is fit using all treatment and control samples except for instance i . Then for each tree in the fitted forest, instance i is placed into its appropriate leaf node in the tree, and the difference between the treated and control outcomes within that node is used as the ITE estimate of that tree. The final estimate is the average estimate of each tree. Let $L(\mathbf{x})$ be the set of instances in the leaf node to which instance i is assigned, $L_1(\mathbf{x}) \subseteq L(\mathbf{x})$ be the subset of treatment samples, and $L_0(\mathbf{x}) \subseteq L(\mathbf{x})$ be the subset of control samples. Then the estimated causal effect of each tree is:

$$\hat{\tau}_{CSF}(\mathbf{x}) = \frac{1}{|L_1(\mathbf{x})|} \sum_{i \in L_1(\mathbf{x})} Y_i - \frac{1}{|L_0(\mathbf{x})|} \sum_{i \in L_0(\mathbf{x})} Y_i \quad (9)$$

In LSE estimation with CSF, after representing each sentence with binary bag-of-words features, we fit a causal forest model to estimate LSE by aggregating estimations from all trees. For our running example, we fit a causal forest using all sentences containing either *shops* or *boutiques*, excluding “*There are plenty of shops nearby*” and then estimate LSE for the sentence by aggregating estimations from all trees, where estimation by each tree is calculated by taking difference between average label values for treatment samples and control samples inside the leaf where “*There are plenty of shops nearby*” belongs to.

4.5 Causal Perception Classifier

The advantages of the approaches above is that they do not require any randomized control trials to collect human perception judgments of lexical substitutions. However, in some situations it may be feasible to perform a small number of RCTs to get reliable LSE estimates for a limited number of sentences. For example, as detailed in §7.2, we can show

subjects two versions of the same sentence, one with w_1 and one with w_2 , and elicit perception judgments. We can then aggregate these into LSE estimates. This results in a set of tuples (w_1, w_2, s, τ) , where τ is the LSE produced by the randomized control trial. In this section, we develop an approach to fit a classifier on such data, then use it to produce LSE estimates for new sentences.

Our approach is to first implement generic, the non-lexicalized features of each (w_1, w_2, s, τ) , then to fit a binary classifier to predict whether a new tuple (w'_1, w'_2, s') has a positive effect on perception or not. This approach is inspired by recent work that frames causal inference as a classification task (Lopez-Paz et al. 2015).

For each training tuple (w_1, w_2, s, τ) , we compute three straightforward features inspired by the intuition of the ITE methods described above. Each feature requires a sentence classifier trained on the class labels (e.g., gender or neighborhood desirability). In our experiments below, we use a logistic regression classifier trained on bag-of-words features.

1. Context probability: The motivation for this feature is that we expect the context in which a word appears to influence its LSE. For example, if a sentence has many indicators that the author is male, then changing a single word may have little effect. In contrast, adding a gender-indicative term to a sentence that otherwise has gender-neutral terms may alter the perception more significantly. To capture this notion, this feature is the posterior probability of the positive class produced by the sentence classifier, using the bag-of-words representation of s after removing word w_1 .

2. Control word probability: This feature is the coefficient for the control word w_1 according to the sentence classifier. The intuition is that if the control word is very indicative of the negative class, then modifying it may alter the perception toward the positive class.

3. Treatment word probability: This feature is the coefficient for the treatment word w_2 according to the sentence classifier. The intuition is that if the treatment word is very indicative of the positive class, then modifying the control word to the treatment word may alter the perception toward the positive class.

We fit a binary classifier using these three features. To convert this into a binary decision problem, we label all tuples where $\tau > 0.5$ as positive examples, and the rest as negative.² To compute the LSE estimate for a new tuple (w_1, w_2, s) , we use the posterior probability of the positive class according to this classifier. See detailed analysis in §8.

5 Data

This section provides a brief description of experimental datasets (Yelp, Twitter, and Airbnb) for LSE estimation.

A key benefit of our first class of approaches is that it does not require data annotated with human perceptions. Instead, it only requires objective annotations. For example, annotations may indicate the self-reported gender of an author, or an objective measure of the quality of a neighborhood, but we do not require annotations of user perceptions of text in

²We use a 1-5 scale in our RCTs, so a treatment effect greater than 0.5 is likely to be significant.

order to produce LSE estimates. While perception and reality are not equivalent, prior work (e.g., in gender perception from text (Flekova et al. 2016)) have found them to be highly correlated. Our results below comparing with human perception measures also support this notion.

Neighborhood Desirability in Airbnb: Airbnb is an online marketplace for short-term rentals, and neighborhood safety is one important factor of desirability that could influence potential guest’s decision. Thus, we use crime rate as proxy of neighborhood desirability. We collect neighborhood descriptions³ from hosts in 1,259 neighborhoods across 16 US cities and collect FBI crime statistics⁴ of each city and crime rate of each neighborhood.⁵ If a neighborhood has a lower crime rate than its city, we label this neighborhood as desirable; otherwise, undesirable. We get 81,767 neighborhood descriptions from hosts in desirable neighborhoods and 17,853 from undesirable neighborhoods.

Gender in Twitter Message and Yelp Reviews: We choose Twitter and Yelp as representative of different social media writing styles to investigate lexical substitution effect on gender perception. First, we use datasets of tweets and Yelp reviews from (Reddy and Knight 2016), where tweets are geo-located in the US and Yelp reviews are originally derived from the Yelp Dataset Challenge released in 2016.⁶ Users in both datasets are pre-annotated with male and female genders. In our sample, we have 47,298 female users with 47,297 male users for Twitter dataset, and 21,650 female users with 21,649 male users for Yelp dataset. Please see Appendix for more details.

6 Generating Candidate Substitutions

Given the combinatorics of generating all possible tuples (w_1, w_2, s) for LSE estimation, we implement several filters to focus our estimates on promising tuples. We summarize these below (see the Appendix for more details):

1. Either w_1 or w_2 must be moderately correlated with the class label (e.g., gender or neighborhood desirability). We implement this by fitting a logistic regression classifier on the labeled data and retaining words whose coefficient has magnitude greater than 0.5.
2. To ensure semantic substitutability, w_2 must be a paraphrase of w_1 , according to the Paraphrase Database (PPDB 2.0) (Pavlick et al. 2015b).
3. To ensure syntactic substitutability, w_1 and w_2 must have the same part-of-speech tag, as determined by the most frequently occurring tag for that word in the dataset.
4. To ensure substitutability for a specific sentence, we require that the n -grams produced by swapping w_1 with w_2 occur with sufficient frequency in the corpus.

After pruning, for Airbnb we obtained 1,678 substitutable word pairs spanning 224,603 sentences from desirable neighborhoods and 49,866 from undesirable neighborhoods; for Twitter we get 1,876 substitutable word pairs

³insideairbnb.com

⁴<https://ucr.fbi.gov/crime-in-the-u.s/2016>

⁵<http://www.areavibes.com/>

⁶https://www.yelp.com/dataset_challenge

spanning 583,982 female sentences and 441,562 male sentences; for Yelp we get 1,648 word pairs spanning 582,792 female sentences and 492,893 male sentences.

7 Experimental Settings

We first carry out experiments to calculate LSE using four estimators and then conduct Randomized Control Trails with Amazon Mechanical Turk (AMT) workers to get human perceived LSE. Next, we fit an out-of-domain causal perception classifier to distinguish LSE directions. Lastly, we evaluate the performance of each method by comparing with human reported values on each dataset separately.

7.1 Calculating LSE Estimates

For experiments with four estimators, we do parameter tuning and algorithm implementation separately. For parameter tuning, we apply the corresponding classification models and do grid search with 5-fold cross validation. For algorithm implementation, we use tuned parameters for each model and follow procedures introduced in §4.

For KNN, we use `KNeighborsClassifier` in `scikit-learn` (Pedregosa et al. 2011) for parameter tuning and then select $k = 30$ for estimator implementation. For VT-RF and CF-RF, we use `RandomForestClassifier` (`scikit-learn`) for parameter tuning and apply the following values in corresponding estimators: $n_estimators = 200$, $max_features = 'log2'$, $min_samples_leaf = 10$, $oob_score = True$. For `CausalForest`, we use the authors’ implementation⁷ and experiment with $n_estimators = 200$ and default values for other parameters as suggested by Wager and Athey (2017).

For the causal perception classifier, our goal is to determine whether the classifier can generalize across domains. Thus, we train the classifier on two datasets and test on the third. We use `scikit-learn`’s logistic regression classifier with the default parameters. To compare this classifier with the results of RCTs, we use the posterior probability of the positive class as the estimated treatment effect, and compute the correlation with RCT estimates.

7.2 Human-derived LSE Estimates

In order to evaluate the methods, and to train the causal perception classifier, we conducted randomized control trails (RCTs) to directly measure how a specific lexical substitution affects reported perceptions. We do so by eliciting perception judgments from AMT workers.

As it would be impractical to conduct AMT for every tuple (w_1, w_2, s) , we instead aim to validate a diverse sample of word substitutions rated highly by at least one of the four LSE estimators. For each dataset, we select the top 10 word substitutions that get the highest LSE according to each estimator. For every selected word substitution (w_1, w_2) , we sample three control sentences (sentences containing w_1) with maximum, minimum and median estimated LSE and generate three corresponding treatment sentences by substituting w_1 to w_2 for each control sentence. Thus, we get

Increase desirability	Increase male perception
store → boutique	gay → homo
famous → grand	yummy → tasty
famous → renowned	happiness → joy
rapidly → quickly	fabulous → impressive
nice → gorgeous	bed → crib
amazing → incredible	amazing → impressive
events → festivals	boyfriends → buddies
cheap → inexpensive	purse → wallet
various → several	precious → valuable
yummy → delicious	sweetheart → girlfriend

Table 1: Samples of substitution words with high LSE

120 control sentences and 120 treatment sentences for each dataset. We divide these sentences into batches of size 10; every batch is rated by 10 different AMT workers. The workers are asked to rate each sentence according to its likely perception of an attribute (on a scale from 1 to 5) (e.g., the neighborhood desirability of an Airbnb description sentence, or the gender of author for Twitter and Yelp sentence). Please see Appendix for details on the annotation guidelines.

For example, for a tuple $(boyfriend, buddy, \text{“My boyfriend is super picky”})$, we have 10 different workers rate the likely gender of the author for *“My boyfriend is super picky”*, then have 10 different workers rate the sentence *“My buddy is super picky”*. The difference in median rating between the second and first sentence is the human perceived effect of changing the word *boyfriend* to *buddy* in this sentence.

Overall, we recruit 720 different AMT workers, 240 for each dataset, and received 237 valid responses for Yelp, 235 for Twitter, and 215 for Airbnb. We compute the Pearson correlation between every two workers who rate the same batch as a measure of inter-annotator agreement as well as the difficulty of LSE tasks for each dataset. These agreement measures, shown in Table 2, suggest that the annotators have moderate agreement (.51-.58 correlation) in line with prior work (Preotiuc-Pietro, Xu, and Ungar 2016). Furthermore, these measures indicate that the Airbnb task is more difficult for humans, which is also expected given that neighborhood desirability is a more subjective concept than gender.

8 Results and Discussion

In this section, we first show a list of substitution words with large LSE estimates, and then provide quantitative and qualitative analysis for different LSE methods.

8.1 Substitution Words with Large LSE

Table 1 shows a sample of 10 substitution words that have large LSE estimates with respect to desirability or gender, based on the automated methods. For example, replacing *shop* with *boutique* increases the perceived desirability of a neighborhood across many sentences. A sentence using the word *tasty* is perceived as more likely to be written by a male than one using *yummy*, and the word *sweetheart* is more often being used by females than *girlfriend*.

⁷<https://github.com/swager/grf>

	Yelp	Twitter	Airbnb
Agreements-pearson	0.557	0.576	0.513
KNN	0.474	0.291	0.076
VT-RF	0.747	0.333	0.049
CF-RF	0.680	0.279	0.109
CSF	0.645	0.338	0.096
Causal perception classifier	0.783	0.21	0.139

Table 2: Inter-annotator agreement and Pearson correlation between algorithmically estimated LSE and AMT judgment

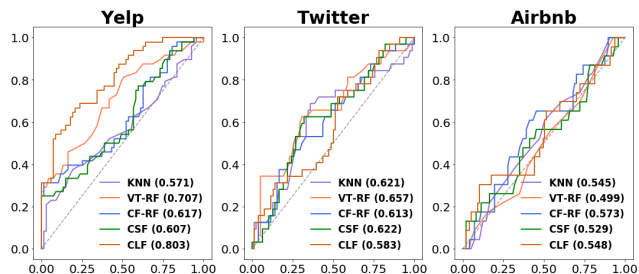


Figure 1: ROC curve for classifying sentences according to AMT perception with estimated LSE as confidence score. (CLF: causal perception classifier). Best viewed in color.

8.2 Comparison with RCT Results

To evaluate the performance of LSE estimators, we first compare algorithmically derived LSE with human derived LSE from AMT. Each tuple (w_1, w_2, s) has both an algorithmically estimated LSE $\hat{\tau}$ by each estimator as well as a human derived LSE τ from AMT workers. For 687 annotated tuples, we calculate the Pearson correlation between algorithmic and human-derived LSE estimates. Table 2 shows the results.⁸ Additionally, Figure 1 plots ROC curves for classifying sentences as having positive or negative treatment effect, using the LSE estimates as confidence scores for sorting instances.

From these results, we can see that LSE estimators are well aligned with human perception measures, which suggests the suitable proxy of algorithmic estimators with perception measure. There is also considerable variation across datasets, with Yelp having the most agreement and Airbnb the least. Yelp has the most formal writing style among the three datasets, so tree-based estimators (CF-RF, VT-RF, CSF) have competitive performance with humans. Twitter is challenging due to grammatical errors and incomplete sentences. Airbnb has less formal writing style compared with Yelp and contains long sentences with proper nouns (e.g., city names, street names and so on) that lead to the lowest correlation and inter-annotator agreement, suggesting that the more subjective the perceptual attribute is, the lower both human agreement and algorithmic accuracy will be.⁹

⁸Human agreement and algorithmic correlations are calculated differently, so the scores may be in slightly different scales.

⁹Using relative crime rates as a proxy for desirability of Airbnb hosts is a possible limitation.

	Yelp	Twitter	Airbnb
context pr	-0.348	-0.829	-0.528
control word pr	-0.141	-0.514	-0.367
treatment word pr	0.189	0.401	0.344

Table 3: Logistic regression coefficients for the features of the causal perception classifier

Additionally, we observe that the causal perception classifier outperforms the four other LSE estimators for two of the three datasets. Table 3 shows coefficient values for the classifier when fit on each dataset separately. These coefficients support the notion that certain aspects of LSE are generalizable across domains — in all three datasets, the sign and relative order of the coefficients are the same. Furthermore, the coefficients support the intuition as to what instances have large, positive effect sizes: tuples (w_1, w_2, s) where w_1 is associated with the negative class (control word probability), where w_2 is associated with the positive class (treatment word probability), and where the context is associated with the negative class (context probability).

Finally, we perform an error analysis to identify word pairs for which the sentence context has a meaningful impact on perception estimates. For example, changing the word *boyfriend* to *buddy* in the sentence “Monday nights are a night of bonding for me and my boyfriend” is correctly estimated to have a larger effect on gender perception than in the sentence “If you ask me to hang out with you and your boyfriend I will ... decline.” The reason is that the use of the possessive pronoun “my” reveals more about the possible gender of the author than the pronoun “your.” We found similar results on Airbnb for the words *cute* and *attractive* — this change improves perceived desirability more when describing the apartment rather than the owner.

9 Conclusion

This paper quantifies the causal effect of lexical change on perception of a specific sentence by adapting concepts from ITE estimation to LSE estimation. We carry out experiments with four estimators (KNN, VT-RF, CF-RF, and CSF) to algorithmically estimate LSE using datasets from three domains (Airbnb, Twitter and Yelp). Additionally, we select a diverse sample to conduct randomized control trials with AMT and fit causal perception classifiers with domain generalizable features. Experiments comparing Pearson correlation show that causal perception classifiers and algorithmically estimated LSE align well with results reported by AMT, which suggests the possibility of applying LSE methods to customize content to perception goals as well as understand self-presentation strategies in online platforms.

Acknowledgments

This research was funded in part by the National Science Foundation under grants #IIS-1526674 and #IIS-1618244.

References

- Afroz, S.; Brennan, M.; and Greenstadt, R. 2012. Detecting hoaxes, frauds, and deception in writing style online. In *Security and Privacy, IEEE Symposium on*, 461–475. IEEE.
- Austin, P. C. 2008. A critical appraisal of propensity-score matching in the medical literature between 1996 and 2003. *Statistics in medicine* 27(12):2037–2049.
- Berger, J., and Milkman, K. L. 2012. What makes online content viral? *Journal of Marketing Research* 49(2):192–205.
- Breiman, L. 2001. Random forests. *Machine Learning* 45(1):5–32.
- Crump, R. K.; Hotz, V. J.; Imbens, G. W.; and Mitnik, O. A. 2008. Nonparametric tests for treatment effect heterogeneity. *The Review of Economics and Statistics* 90(3):389–405.
- Danescu-Niculescu-Mizil, C.; Cheng, J.; Kleinberg, J.; and Lee, L. 2012. You had me at hello: How phrasing affects memorability. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, 892–901. Association for Computational Linguistics.
- Flekova, L.; Carpenter, J.; Giorgi, S.; Ungar, L.; and Preotiuc-Pietro, D. 2016. Analyzing biases in human perception of user age and gender from text. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, volume 1, 843–854.
- Foster, J. C.; Taylor, J. M.; and Ruberg, S. J. 2011. Subgroup identification from randomized clinical trial data. *Statistics in Medicine* 30(24):2867–2880.
- Hagen, M.; Potthast, M.; and Stein, B. 2017. Overview of the author obfuscation task at pan 2017: safety evaluation revisited. *Working Notes Papers of the CLEF* 33–64.
- Hovland, C.; Janis, I.; and Kelley, H. 1953. *Communication and persuasion: psychological studies of opinion change*. Greenwood Press.
- Lee, M.-J. 2008. Non parametric tests for distributional treatment effect for randomly censored responses. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 71(1):243–264.
- Lopez-Paz, D.; Muandet, K.; Schölkopf, B.; and Tolstikhin, I. 2015. Towards a learning theory of cause-effect inference. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *JMLR Workshop and Conference Proceedings*, 1452–1461. JMLR.
- Lu, M.; Sadiq, S.; Feaster, D. J.; and Ishwaran, H. 2018. Estimating individual treatment effect in observational data using random forest methods. *Journal of Computational and Graphical Statistics* 27(1):209–219.
- Ludwig, S.; De Ruyter, K.; Friedman, M.; Brügger, E. C.; Wetzels, M.; and Pfann, G. 2013. More than words: The influence of affective content and linguistic style matches in online reviews on conversion rates. *Journal of Marketing* 77(1):87–103.
- Nguyen, T. T. D. T.; Garncarz, T.; Ng, F.; Dabbish, L. A.; and Dow, S. P. 2017. Fruitful feedback: Positive affective language and source anonymity improve critique reception and work outcomes. In *Proceedings of the 2017 ACM Conference on Computer Supported Cooperative Work and Social Computing*, 1024–1034. ACM.
- Packard, G., and Berger, J. 2017. How language shapes word of mouth’s impact. *Journal of Marketing Research* 54(4):572–588.
- Pavlick, E.; Bos, J.; Nissim, M.; Beller, C.; Van Durme, B.; and Callison-Burch, C. 2015a. Adding semantics to data-driven paraphrasing. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, 1512–1522. Association for Computational Linguistics.
- Pavlick, E.; Rastogi, P.; Ganitkevitch, J.; Van Durme, B.; and Callison-Burch, C. 2015b. Ppdb 2.0: Better paraphrase ranking, fine-grained entailment relations, word embeddings, and style classification. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, 425–430. Association for Computational Linguistics.
- Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. 2011. Scikit-learn: Machine learning in python. *Journal of machine learning research* 12(Oct):2825–2830.
- Preotiuc-Pietro, D.; Guntuku, S. C.; and Ungar, L. 2017. Controlling human perception of basic user traits. In *Proceedings of the 2017 conference on Empirical Methods in Natural Language Processing*, 2335–2341.
- Preotiuc-Pietro, D.; Xu, W.; and Ungar, L. H. 2016. Discovering user attribute stylistic differences via paraphrasing. In *AAAI*, 3030–3037.
- Reddy, S., and Knight, K. 2016. Obfuscating gender in social media writing. In *Proceedings of the First Workshop on NLP and Computational Social Science*, 17–26.
- Riley, T. J., and Luippold, B. L. 2015. Managing investors’ perception through strategic word choices in financial narratives. *Journal of Corporate Accounting & Finance* 26(5):57–62.
- Rubin, D. B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology* 66(5):688.
- Tan, C.; Lee, L.; and Pang, B. 2014. The effect of wording on message propagation: Topic- and author-controlled natural experiments on twitter. In *Proceedings of ACL*.
- Thibodeau, P. H., and Boroditsky, L. 2013. Natural language metaphors covertly influence reasoning. *PLoS one* 8(1):e52961.
- Wager, S., and Athey, S. 2017. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*.
- Willke, R. J.; Zheng, Z.; Subedi, P.; Althin, R.; and Mullins, C. D. 2012. From concepts, theory, and evidence of heterogeneity of treatment effects to methodological approaches: a primer. *BMC Medical Research Methodology* 12(1):185.

A Appendix A: Additional Results

We provide supplemental information and detailed analysis in this section.

A.1 Details for Datasets

Airbnb We collect neighborhood descriptions from hosts in 1,259 neighborhoods across 16 US cities from insideairbnb.com by May 2017. Table 4 shows the 16 cities and corresponding number of neighborhoods we collect in each city.

City	Number of Neighborhoods
LA	248
NY	219
Oakland	108
SanDiego	96
Portland	92
Seattle	87
Denver	74
Chicago	74
NewOrleans	69
Austin	43
WDC	39
SanFrancisco	37
Nashville	35
Boston	25
Asheville	8
SantaCruz	5

Table 4: Cities and number of neighborhoods in each city

Neighborhood Desirability As a subjective concept, desirability of a rental could be measured by multiple factors such as safety, convenience, surroundings, traffic and so on. In this paper, we aim to get an objective measure that could be applied to rentals anywhere and since we only consider Airbnb rentals inside USA, where safety is a very important factor that could influence potential guest’s decision, so we decide to use relative crime rate as proxy of neighborhood desirability. However, we acknowledge the limitation of making this assumption. A rental that is attractive to one person who prefer safety might not be attractive to another who prefer location.

We collect crime rate of cities and neighborhoods separately from two sources. For crime rate of cities, we collect from FBI crime statistics¹⁰. For crime rate of each neighborhood, we collect from areavibes¹¹. Considering that crime rate varies from city to city, it is unfair to directly compare neighborhoods in different cities, we make comparisons inside each city by comparing the relative crime rate of a neighborhood with the city it locates as our labeling criteria. We conduct the labeling process as follows:

- Label a neighborhood: if a neighborhood has lower crime rate than the city it locates, we label this neighborhood as desirable; otherwise, undesirable.

¹⁰<https://www.freep.com/story/news/2017/09/25/database-2016-fbi-crime-statistics-u-s-city/701445001/>

¹¹<http://www.areavibes.com/>

- Label a host: we assign the same label for hosts located in one neighborhood and get 81,767 neighborhood descriptions from hosts in desirable neighborhoods and 17,853 from undesirable neighborhoods. We observe the data imbalance which might be due to the fact that low-crime areas are more desirable to potential guests, so more Airbnb rentals are listed in low-crime areas than in high-crime areas.
- Label a sentence: we label each neighborhood description sentence by the label of that neighborhood, which means all desirable neighborhood description sentences are labeled as desirable, otherwise undesirable.

Twitter and Yelp We use tweets and Yelp reviews from datasets introduced in (Reddy and Knight 2016). According to (Reddy and Knight 2016), tweets are collected in July 2013 and only consider those geolocated in US; the corpus of Yelp reviews is obtained from 2016 Yelp Dataset Challenge¹².

The two datasets are annotated with two genders: male and female, which are inferred by mapping users’ first names with Social Security Administration list of baby names from 1990¹³. While male and female are suggested as accurate reflection of social media users’ genders, we consider non-binary gender labels as an important area of future work.

After processing by removing users with ambiguous names, dropping non-English and highly gendered texts, they get 432,983 user corpus for Yelp and 945,951 for Twitter (Reddy and Knight 2016). Sampling from their datasets, we get Twitter corpus from 47,298 female users and 47,297 male users, and Yelp corpus from 21,650 female users and 21,649 male users.

Please refer to (Reddy and Knight 2016) for more details about Twitter and Yelp datasets.

A.2 Identify Qualified Lexical Substitutions

To generate tuples $(w_1, w_2, sentence)$ for LSE estimation tasks, we first search for substitutable word pairs (w_1, w_2) and then select sentences that are qualified for substituting w_1 to w_2 .

Select Representative Words Considering the large number of possible lexical substitutions, we first apply several criteria to select the most representative words and then match them with the most appropriate substitutions. To explore the subtle effect of a single word change on perceived perception of the corresponding sentence, we first select words that are representative of attributes we are interested in and thus substituting them might cause effects large enough to be captured. For example, given a sentence: “*I had lunch with my boyfriend*” written by female, *boyfriend* is the most representative words with regard to gender of female, substituting *boyfriend* to *girlfriend* will change the perceived gender of the author from female to male, while substituting the word *had* to *took* does not change the perceived perception. To select representative words, we fit a

¹²https://www.yelp.com/dataset_challenge

¹³<https://www.ssa.gov/oact/babynames/limits.html>

	Airbnb	Twitter / Yelp
Q&A	Rate the desirability of a short-term apartment rental based on a single sentence.	Rate how likely you think this tweet / Yelp review sentence is written by male or female.
5	Very desirable	Very likely male
4	Somewhat desirable	Somewhat likely male
3	Neither desirable nor undesirable	Neutral, neither male nor female
2	Somewhat undesirable	Somewhat likely female
1	Very undesirable	Very likely female

Table 5: Amazon Mechanical Turk annotation guidelines

binary Logistic Regression classifier for each dataset separately.

For Airbnb dataset, we fit a classifier with 81,767 desirable and 17,853 undesirable neighborhood descriptions. Considering that description texts contain lots of proper nouns like street names, famous place names, neighborhood names and city names, we limit the vocabulary to common words that appear at least 8 times in 6 cities and thus eliminating classifier bias towards proper nouns. By doing so, we get 1,549 common words as representative words of desirable and undesirable classes.

For Twitter and Yelp datasets, after marking proper nouns with NLTK toolkit¹⁴, we fit a binary classifier for Twitter with tweets from 47,298 female users and 47,297 male users. And a classifier for Yelp with reviews from 21,650 female users and 21,649 male users. Using coefficient thresholds greater than 0.5 or smaller than -0.5, we select 4,087 gender representative words for Twitter and 2,264 for Yelp.

After selecting representative words, we search for semantically and syntactically qualified substitutions for them.

Semantically Qualified Substitutions (Reddy and Knight 2016) apply word2vec extensions of Yelp reviews and tweets parsed with CoreNLP and TweetNLP to capture semantically similar words, and (Preotiuc-Pietro, Xu, and Ungar 2016) use Paraphrase Database (PPDB) to get stylistic paraphrases with equivalence probability greater than 0.2. In our case, we have three corpus with different writing styles and our goal is to find single word substitutions that express the same meaning, so we choose PPDB as our source to get paraphrases in this paper and will consider word2vec extensions of Airbnb corpus in future work. PPDB((Pavlick et al. 2015a)) is a collection high precision paraphrases extracted from bilingual parallel corpora with each paraphrase be assigned with probability and similarity scores according to Google ngrams and Gigaword corpus, and later extended with equivalent scores that interpret semantic relationship between paraphrase pairs. We search for paraphrase pairs with equivalence probability of at least 0.15 ((Preotiuc-Pietro, Xu, and Ungar 2016) use 0.2, we decide to use 0.15 as a relative loose criteria).

Syntactically Qualified Substitutions Despite of checking semantics of substitution words, we need to make sure the substitutions are also syntactically qualified. For example, substitutable words should have same singular or plural

forms. To do so, we first do POS tagging¹⁵ for all sentences in three corpus and store the annotated POS tags of each word, and then check the most common POS tag of each paraphrase pair and only retain paraphrase pairs that have the same most common POS tags.

After limiting substitutions of representative words to semantically and syntactically suitable ones, we search for sentences that are qualified for each specific word substitution.

Check Word Substitutability in Specific Sentences We first build a bi-gram vocabulary using three datasets. Then, for each pair of substitution words (w_1, w_2) , we search for sentences containing w_1 and check for every sentence that if substituting w_1 to w_2 produces valid bi-grams by looking up the bi-gram vocabulary. For example, to check the substitutability of $(perced, drilled)$ in “*I’m having my ears perced on Saturday*”, we decide the grammatical correctness of the sentence after substitution “*I’m having my ears drilled on Saturday*” by checking if “*ears drilled*” and “*drilled on*” exist in our bi-gram vocabulary. If yes, we will keep the current sentence as a qualified sentence for this substitution, otherwise, remove the sentence.

Overall, after pruning with the above criteria, we obtained 1,678 substitutable word pairs spanning 224,603 sentences from desirable neighborhoods and 49,866 from undesirable neighborhoods; and 1,876 substitutable word pairs spanning 583,982 female sentences and 441,562 male sentences for Twitter dataset; and 1,648 word pairs spanning 582,792 female sentences and 492,893 male sentences for Yelp dataset.

A.3 Crowd-sourcing Experiments with Amazon Mechanical Turk

We take a tuple $(w_1, w_2, sentence)$ as the unit of analysis in LSE tasks. Despite of algorithmically calculate how much does substituting w_1 to w_2 for the *sentence* affects its perceived perception, we conduct Randomized Control Trails to directly measure LSE by eliciting judgments from Amazon Mechanical Turk (AMT) workers. Detailed procedures are as follows:

- **Select word pairs with highest LSE** Among all substitution word pairs, we first select those rated highly by at least one of the four LSE estimators (KNN, VT-RF, CT-RF, CSF). Specifically, for each dataset, we get top-10

¹⁴<https://www.nltk.org/>

¹⁵We use NLTK (<http://nltk.sourceforge.net/>) for POS tagging.

word substitutions according to each of the four estimators. If a substitution word pair is rated as top-10 with more than one estimators, then we only keep this word pair for the estimator that gives the highest rank (e.g., for a substitution word pair (w_1, w_2) , if KNN estimator rank it as the second and VT-RF estimator ranks it as the fifth, then we keep (w_1, w_2) for KNN estimator). Thus, we get 10 substitution word pairs for each of the four estimators.

- **Select sentences with maximum, minimum and median LSE for each word pair** For each word substitution (w_1, w_2) , we rank all control sentences (e.g., sentences containing w_1) according to LSE calculated by the corresponding estimator and sample three sentences with maximum, minimum and median LSE. Meanwhile, we generate corresponding treatment sentences using the given substitution word (w_2) . Thus, we get 120 control sentences and 120 treatment sentences for each dataset.
- **Create AMT tasks** For each dataset, we divide 120 control sentences into 12 batches with each batch has 10 different sentences, and the same process for 120 treatment sentences. We take each batch as a HIT task in AMT, and for each HIT task, we recruit 10 different workers and ask them to pick a scale (ranges from 1 to 5) for every sentence according to its likely perception of an attribute. Table 5 shows the annotation guidelines for three datasets.
- **Quality control of AMT tasks** To eliminate possible biases, we limit that each worker only have access to one batch of either control or treatment sentences. If a worker rates a batch of control sentences, then he won't be able to see the corresponding treatment sentences, so that his decision is not affected by knowing which word is being substituted. For quality control, we require workers to be graduates of U.S. high schools, and we include attentiveness checks using manually created "dummy" sentences. For example, a "dummy" sentence for gender perception, "I am the son of my father", should be rated as written by a male. We remove responses from workers who provide incorrect answers for dummy questions.

A.4 Experiments with LSE Estimators

We first conduct parameter tuning to select the most suitable parameters for each estimator and then implement four estimators following procedures introduced in §4.

Parameter Tuning As we are estimating LSE on sentence level, we do parameter tuning with all labeled sentences of each dataset. Parameters are tuned for the classification task, but not for the treatment effect estimation task (none of the KNN/VT-RF/CF-RF/CSF methods were tuned using the labeled AMT data, so we can measure effectiveness without access to such expensive data).

- **Feature Representation** We try both bag-of-words and tf-idf feature representation techniques for each method.
- **KNN tuning** We use scikit-learn implementation of KNeighborsClassifier and do grid search for $n_neighbors$ (since we only need the number of neighbors in KNN estimator implementation, so we don't consider other param-

Yelp	Label
My wife likes this place.	Male
I like coming here with my fraternity brothers.	Male
My brother and I come here for guys night out.	Male
My husband likes this place.	Female
I like coming here with my sorority sisters.	Female
My sister and I come here for girl's night out.	Female
Twitter	
I love playing football and video games.	Male
My wife is waiting on me.	Male
I am my father's son.	Male
I love getting a pedicure at girls night out.	Female
My husband says I smile too much.	Female
I am my mom's daughter.	Female
Airbnb	
This is by far the best neighborhood in the city.	Desirable
This neighborhood is amazing in every way.	Desirable
What a world-class neighborhood this is!	Desirable
This neighborhood is not so great.	Undesirable
Yes, there is a lot of crime in this neighborhood.	Undesirable
Lots of shootings in this neighborhood.	Undesirable

Table 6: Dummy sentences for Yelp, Twitter and Airbnb

eters) and get the best 5fold cross validation score with $n_neighbors = 30$.

- **Random-Forest tuning** We use scikit-learn implementation of Random Forest classifier and do grid search for a set of parameters and get the best 5-fold cross validation score with $n_estimators = 200$, $max_features = 'log2'$, $min_samples_leaf = 10$ and $oob_score = True$.

As mentioned in previous context, there exists imbalance between the number of 'desirable' and 'undesirable' descriptions in Airbnb dataset. We considered model variants that deal with class imbalance (e.g., overweighting the minority class), but did not observe significantly different results with such methods.

Estimator Implementation For estimator implementation, we follow the process introduced in §4 and use the best parameters reported by the above tuning process for KNN VT-RF, CF-RF. For Causal Forest, we try $n_estimators = 200$ with default values of other parameters.

A.5 Causal Perception Classifier

We fit two classifiers for this task. First, we fit one classifier for each dataset to get proposed features: posterior probability of a context, coefficient of substitution words, and the number of positively and negatively related words. After representing each tuple $(w_1, w_2, sentence)$ with proposed features, we fit causal perception classifiers only using samples labeled by Amazon Mechanical Turks. Specifically, each causal perception classifier is fitted by using samples of two datasets and making out-of-domain prediction for the third dataset.

Yelp	Twitter	Airbnb
lovely → delightful	gay → homo	store → boutique
cute → attractive	yummy → tasty	famous → grand
helpful → useful	happiness → joy	famous → renowned
fabulous → terrific	fabulous → impressive	rapidly → quickly
gorgeous → outstanding	bed → crib	nice → gorgeous
salesperson → dealer	amazing → impressive	amazing → incredible
belongings → properties	boyfriends → buddies	events → festivals
thorough → meticulous	purse → wallet	cheap → inexpensive
happily → fortunately	precious → valuable	various → several
dirty → shitty	sweetheart → girlfriend	yummy → delicious
Increase male perception or decrease female perception		Increase desirability

Table 7: Substitutable word pairs with large LSE

A.6 Results and Analysis

In this section, we provide both qualitative and quantitative analysis from the following aspects:

- First, we present a sample of substitution words estimated to have large LSE.
- Second, we compare the performance of four LSE estimators.
- Third, we evaluate the agreement of each estimator with human perception RCTs using Amazon Mechanical Turk.
- Fourth, we assess the causal perception classifier and interpret feature importance with experimental findings.
- Finally, we provide a preliminary analysis of how this approach may be used to characterize communication strategies online.

Substitution Words with Large LSE Table 7 shows a sample of substitutable word pairs estimated to have large LSE by at least one estimator.

For Airbnb, the substitution words are reported to increase the perceived desirability of a rental. For example, since *boutique* often related with nice neighborhoods, substituting *shop* to *boutique* helps increase the neighborhood desirability. For Twitter and Yelp, the substitution words are reported to increase male perception or decrease female perception of the author. For example, a sentence using *tasty* is more likely to be written by a male than using *yummy*, and chances are high that *sweetheart* would appear in a female sentence while *girlfriend* in a male sentence.

Additionally, to assess the quality of substitutable word pairs, we select top 20 word pairs with largest LSE reported by each estimator and manually check if these word pairs are both syntactically and semantically qualified substitutions. As indicated by Table 8, we find that KNN estimator is somewhat more likely to assign large LSE for qualified substitutions. Unsuitable word pairs are often generated due to the fact that the paraphrase database (PPDB) was trained on general texts, but the validity of a substitution can depend on domain. For example, *gross* and *overall* are potential paraphrases according to PPDB due to one sense of *gross*, but in the Twitter data *gross* is instead more commonly used as a synonym for *disgusting*. More conservative pruning using

	Yelp	Twitter	Airbnb	Mean
KNN	100%	85%	90%	91.67%
VT-RF	100%	65%	90%	85%
CF-RF	85%	75%	75%	78.33%
CSF	80%	70%	50%	66.67%

Table 8: Fraction of top 20 substitutable word pairs that are judged to be acceptable by manual review

language models trained on the in-domain data may reduce the frequency of such occurrences.

Quantitative Analysis of LSE Estimators In this section, we quantitatively compare the similarities and differences between four estimators. We expect there to be differences between KNN and the forest-based methods, since their underlying classification functions are different: KNN estimator directly search from all training instances to identify k nearest neighbors in control and treatment group. In contrast, VT-RF, CF-RF and CSF are all tree-based methods, which attempt to place instances in the same leaf if they are homogeneous with respect to the covariate vector \mathbf{X} .

To quantitatively compare the performance of four estimators, we first generate the entire ranked list of $(w_1, w_2, sentence)$ tuples according to each estimator and then compute Spearman’s rank correlation for ranked list of every two estimators.

According to results shown in Table 9, we observe that:

- Forest based methods (VT-RF, CF-RF, CSF) perform more similar than KNN.
- Four estimators have less agreement on Airbnb dataset than on Twitter and Yelp, which suggests that estimating LSE on Airbnb is harder, because hosts are incentivized to highlight desirable aspects of the neighborhood.

Then, we calculate the percentage of sentences labeled as negative (refers to undesirable for Airbnb and female for Yelp and Twitter) among top 1000 sentences with large LSE. Results in Table 10 shows that:

- All of the four estimators tend to pick negative instances for large LSE. Since we rank sentences in descending order of estimated LSE, the more number of negative sen-

	Yelp				Twitter				Airbnb			
	KNN	VT-RF	CF-RF	CSF	KNN	VT-RF	CF-RF	CSF	KNN	VT-RF	CF-RF	CSF
KNN	1.0	0.674	0.715	0.655	1.0	0.699	0.729	0.668	1.0	0.469	0.561	0.455
VT-RF		1.0	934	0.945		1.0	0.932	0.935		1.0	0.822	0.773
CF-RF			1.0	0.899			1.0	0.883			1.0	0.733
CSF				1.0				1.0				1.0

Table 9: Spearman correlation between ranked sentences of four estimators

tences ranked in top 1000 the more effective that estimator is.

- CF-RF estimator picks the most negative instances for large LSE.
- VT-RF estimator performs differently with other estimators, and especially for Airbnb dataset. The reason may lie in the fact that we label each description sentences as desirable or undesirable according to relative crime rate of a neighborhood, which means all sentences describing low-crime neighborhoods are labeled as desirable and vice versa. However, sentences describing low-crime neighborhoods are not guaranteed to disclose desirability but will be mislabeled as desirable according to our criteria, and this misleads VT-RF estimator and explains the difference of this estimator.

	Yelp	Twitter	Airbnb
KNN	90.5%	70.5%	86.6%
VT-RF	93.9%	71.3%	64.9%
CF-RF	96.6%	77.3%	87.7%
CSF	95.9%	71%	84.4%

Table 10: Percentage of negative sentences in top 1000 highly ranked instances with respect to LSE

Qualitative Analysis of Estimators To qualitatively assess the performance four estimators, we first show examples to get a better understanding of how do four estimators perform differently in recommending substitutable words for a sentence. As shown in Table 14:

- For Yelp, we pick a sentence labeled as male and find substitutable words to make it more likely a sentence written by female. Four estimators give same recommendations for this sentence.
- For Twitter, we pick a sentence labeled as female and four methods recommend substitutable words to make it more likely a male sentence. E.g., as *boyfriend* is most likely to be used by females while *buddy* by males, substituting *boyfriend* to *buddy* makes the sentence more likely to be perceived as written by male.
- For Airbnb, we pick a neighborhood description sentence labeled as undesirable, and four estimators make recommendations to improve its desirability. CF-RF and CSF agree on recommendations for this sentence.

Additionally, we show an example in table 15 to see how do LSE vary for same word substitution in different sentences. We randomly pick one substitutable word pair in

each dataset, and get its highest and lowest LSE sentence according to CSF estimator.

- For Airbnb, substituting *shop* to *boutique* gives lowest LSE on the sentence that is less immediately associated with rental, because it is “located a mile away”.
- For Twitter, substituting *boyfriend* to *buddy* gives highest treatment effect for the sentence talking about “my boyfriend”, which the word “my” is directly associated with the writer of this sentence, so substituting it to “my buddy” makes a big change on the writer’s gender. But for the lowest treatment sentence, the substitution makes a small change because “your boyfriend” and “your buddy” do not refer to the writer’s gender.

Performance of Causal Perception Classifier Our goal for causal perception classifier was to use a small number of generic features to allow the method to generalize across domains (e.g., we fit a model on Yelp and apply it to Twitter). Despite of results shown in §8, we performed some preliminary experiments with a few other features (e.g., sentence length, part-of-speech, number of support words and conflict words), but did not observe significantly different results.

Performance	Yelp	Twitter	Airbnb
AUC	0.803	0.583	0.548
Precision	0.80	0.70	0.65
Recall	0.69	0.72	0.81
F1	0.63	0.62	0.72

Table 11: Performance of causal perception classifier

Comparing LSE Reported by Estimators with Human Judgments In this section, we evaluate the agreement between LSE estimators with human perception RCTs using Amazon Mechanical Turk. To do this, we first calculate inner-annotator agreement using both pearson and Spearman’s rank and take it as a measure of the difficulty of LSE task with each dataset, and then compute Pearson correlation between LSE reported by four estimators. For the RCTs, we compute human perceived LSE as the difference between median ratings for treatment sentence and control sentence.

Table 2 shows the Pearson correlation between each LSE estimator and AMT reported LSE and Figure 2 shows the extend ROC curve for classifying sentences according to AMT perception with estimated LSE as confidence score.

- LSE estimated by four estimators are well aligned with AMT perceived results, which suggests the suitable proxy of objectives measures we use with perception measure.

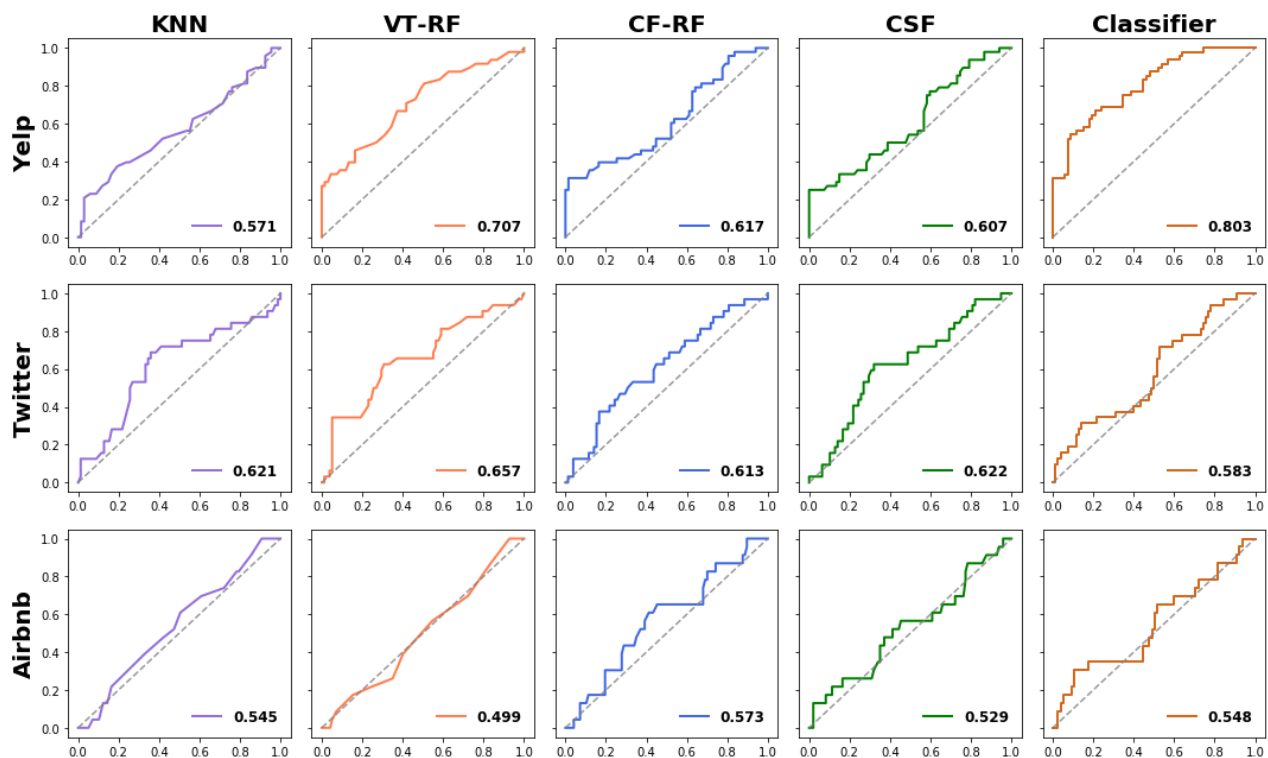


Figure 2: ROC curve for classifying sentences according to AMT perception with estimated LSE as confidence score

Specially treatment effect for Yelp dataset calculated by CF-RF method has the highest correlation 0.57.

- LSE task for Yelp has the highest correlation with AMT perceived results and three tree-based methods (CF-RF, VT-RF, CSF) have competing performance with inter-annotator agreement.
- LSE task for Airbnb has the lowest correlation, and inter-annotator agreement by Pearson and Spearman’s rank also be the lowest, which suggests the difficulty for LSE estimation on Airbnb dataset.
- The more subjective the perceptual attribute is, the lower both human agreement and machine accuracy will be. Since Yelp has the most formal writing style among the three datasets, LSE estimators perform as good as humans. Twitter is challenging due to its informal writing, as compared with Yelp, which contains more grammatically correct, complete sentences. Beyond the fact that desirability is subjective, Airbnb has an informal writing style and contains long sentences with proper nouns (e.g., city names, street names and so on), which decrease the sentence readability for AMT workers. Besides that, since hosts are motivated to attract guests by highlighting positive aspects and roundabout negative aspects, the use of euphemism increases the difficulty of this task: on one hand, it increases the difficulty for human understanding; on the other hand, it misleads LSE estimators as we did not equip LSE algorithms with the ability to identify euphemisms.

- Estimators’ performance with Yelp dataset correlate with humans more than humans correlate with each other. We have two possible explanations for this: First, human agreement is calculated from the average pairwise correlation across 10 AMT workers annotating the same 10 sentences. In contrast, the algorithmic correlations are calculated by comparing the algorithmic scores with the median human scores across 200 or so sentences. Because of this somewhat different calculation, the scores may be in slightly different scales. Second, while we implemented several quality control measures for AMT (see the end of section A.3 in the supplementary material), there are still some outlier workers who reduce the overall agreement number. This in part motivates our use of the median rating when computing the final results.

In addition to correlation, we also evaluate whether the sign of algorithmically estimated LSE agree with AMT perceived LSE. To do so, we code estimated LSE as positive or negative, and compute ROC curves for each estimator shown in Table 13.

	Yelp	Twitter	Airbnb
KNN	0.571	0.621	0.545
VT-RF	0.707	0.657	0.499
CF-RF	0.617	0.613	0.573
CSF	0.607	0.622	0.529
Classifier	0.803	0.583	0.548

Table 13: Area under ROC curve

	Airbnb	Twitter	Yelp
Increase desirability or male perception	closest → best stores → boutiques famous → old plaza → place	okay → good sweatheart → girlfriend purse → wallet precious → rare	gorgeous → super yummy → tasty fabulous → excellent hunt → search
Decrease desirability or male perception	excellent → safe best → hottest gorgeous → great boutiques → stores	ma → mom crib → bed impressive → wonderful buddy → boyfriend	tasty → yummy excellent → cute good → yummy attractive → cute

Table 12: Word substitutions with high LSE used most frequently by authors of the opposite class (e.g., “male” words used by female users, and visa versa.)

Preliminary Analysis using LSE in Online Communication Strategy

In this section, we provide a preliminary analysis of how LSE estimators may be used to characterize communication strategies online. We show potential communication strategies people use for perception management (try to improve positive perception and reduce negative perception, or to change female style to male or vice versa) according to results suggested by current datasets.

To do this, we first select top 20 highest and lowest ranked substitutable word pairs according to each LSE estimator. Then, for the 20 highest ranked word pairs, we sort them according to the frequency of positive treatment words used in negative sentence; for the 20 lowest word pairs, we sort them according to the frequency of negative treatment words used in positive sentence. Table 12 shows a list of highly ranked word pair selected according to each estimator:

- For Airbnb, hosts in undesirable neighborhoods use words *best* instead of *closest* and *boutiques* instead of *shops* more often, which are signs of improving desirable perception. While for hosts in desirable neighborhoods, the estimators suggest them to use words *excellent* instead of *safe* because *safe* reduces positive perception compared with *excellent* (according to LSE recommendations). This makes sense because hosts located in safe neighborhoods would not emphasize safety.
- For gender perception of Twitter and Yelp, LSE estimators recommend that if you want to write sentence like a female, then use *sweatheart* instead of *girlfriend* and use *yummy* instead of *tasty*. Otherwise, if you want to write sentences like a male, use *buddy* instead of *boyfriend* and use *attractive* instead of *cute*. Additionally, LSE estimators recommend to use more emotional words for female sentence than for male.

Yelp (make it more likely a female sentence)	
Original	Very fresh , and <u>tasty</u> herbs and spring rolls as well !
KNN/VT-RF/ CF-RF/CSF	Very fresh , and <u>yummy</u> herbs and spring rolls as well !
Twitter (make it more likely a male sentence)	
Original	Every girl I know is with it and makes <u>nice</u> dinners for their <u>boyfriends</u> while I just order pizza and drink <u>beer</u> with mine #sorrybabe.
KNN/CF-RF	Every girl I know is with it and makes <u>good</u> dinners for their <u>buddies</u> while I just order pizza and drink <u>beer</u> with mine #sorrybabe.
VT-RF/CSF	Every girl I know is with it and makes <u>nice</u> dinners for their <u>buddies</u> while I just order pizza and drink <u>brew</u> with mine #sorrybabe.
Airbnb (increase desirability)	
Original	I don't suggest long walks after dark, but I would <u>definitely</u> not let this neighborhood discourage your stay, it's <u>vibrant</u> , fun and <u>exciting</u> .
KNN	I don't suggest long walks after dark, but I would <u>truly</u> not let this neighborhood discourage your stay, it's <u>dynamic</u> , fun and <u>interesting</u> .
VT-RF	I don't suggest long walks after dark, but I would <u>really</u> not let this neighborhood discourage your stay, it's <u>dynamic</u> , fun and <u>stunning</u> .
CF-RF/CSF	I don't suggest long walks after dark, but I would <u>absolutely</u> not let this neighborhood discourage your stay, it's <u>dynamic</u> , fun and <u>spectacular</u> .

Table 14: Different recommendations of substitution words for one sentence

Yelp (cute → attractive)	
Largest effect	The joint is <u>cute</u> and clean and parking is a breeze.
Smallest effect	Our <u>cute</u> Long Island native , Mary suggested the best things on the menu - even telling us what was off and on from the specials board that would work or not.
Twitter (boyfriend → buddy)	
Largest effect	Monday nights are a night of bonding for me and my <u>boyfriend</u> ! We both LOVE #TeenWolf user user.
Smallest effect	If you ask me to hang out with you and your <u>boyfriend</u> I will look at you like you're stupid then impolitely decline.
Airbnb (store → boutique)	
Largest effect	Check: Andersonville, in particular, has attracted many gay residents (who have re-made the upper reaches of Clark Street into a hot design- <u>store</u> destination).
Smallest effect	Beachwood Village grocery store and coffee <u>shop</u> conveniently located a mile away.

Table 15: Sentences that get the largest and smallest treatment effects for a same word pair