1 DETECTING BEHAVIORAL FAILURES IN EMERGING ELECTRIC VEHICLE

- 2 INFRASTRUCTURE USING SUPERVISED TEXT CLASSIFICATION ALGORITHMS 3
- 4
- 5
- 6 Sooji Ha
- 7 School of Civil & Environmental Engineering and School of Computational Science &
- 8 Engineering
- 9 Georgia Institute of Technology, Atlanta, GA, 30332
- 10 Email: s.ha@gatech.edu
- 11

12 Daniel J. Marchetto

- 13 School of Public Policy
- 14 Georgia Institute of Technology, Atlanta, GA, 30332
- 15 Email: daniel.marchetto@gatech.edu
- 16

17 Mary Elizabeth Burke

- 18 School of Public Policy
- 19 Georgia Institute of Technology, Atlanta, GA, 30332
- 20 Email: mburke38@gatech.edu
- 21

22 Omar Isaac Asensio*

- 23 School of Public Policy and Institute for Data Engineering & Science (IDEaS)
- 24 Georgia Institute of Technology, Atlanta, GA, 30332
- 25 Email: asensio@gatech.edu
- 26 ORCID Number: 0000-0003-2143-5022
- 27 * Corresponding Author
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36

1 ABSTRACT

- 2 There is a growing interest in applying computational tools to the automatic discovery of social
- 3 and economic behavior. For example, with decisions involving resource allocation related to pub-
- 4 lic infrastructure, the ability to predict failures can allow for more efficient policy responses. In
- 5 this paper, we use social data from a popular electric vehicle (EV) driver app to characterize the
- 6 emerging EV charging station infrastructure. We introduce a typology of EV charging experiences
- 7 collected from user reviews and deploy text classification algorithms, including convolutional neu-
- 8 ral networks (CNN), to automatically learn about potential failures. We use machine learning 9 techniques as a pre-processing tool for econometric analyses on the quality of service delivery.
- 10 After classifying the reviews into 9 main user topics and 34 subtopics, we find that the dominant
- 11 issues in EV charging relate to station functionality and availability, which drive negative consumer
- 12 experience. Contrary to the public discourse about EVs, range anxiety was not of large concern to
- 13 existing EV drivers. Based on our findings, we move towards automated identification of failures
- 14 in public charging infrastructure that can significantly reduce research evaluation costs through
- 15 relatively simple computational solutions.
- 16
- 17 Keywords: electric vehicles, consumer behavior, convolutional neural networks, natural language
- 18 processing, mobile data

1 INTRODUCTION

2 The transportation sector has become a dominant source of CO₂ emissions in the United States 3 (1). In the last few years, there has been a growing attention on vehicle electrification as a strategy to reduce mobile source emissions with positive spillovers in air quality benefits (2). For example, 4 in order to accelerate electric vehicle (EV) purchases, a majority of U.S. states are offering some 5 type of financial incentive to complement federal tax credits, including rebates, tax exemptions and 6 other incentives (3-5). An important complementarity to EV adoption is the availability of public 7 charging infrastructure. It is now estimated that global investment in EV charging infrastructure 8 9 from public and private sources will reach \$80 billion USD by 2025 (6). In the US, this investment 10 growth also marks an expected transition in policy support to a focus on charging infrastructure. However, currently there is no easy way to evaluate the needs of drivers or assess perfor-11 mance. This is because the infrastructure upgrades needed to be able to monitor electric consump-12 tion and use in individual charging stations are at early stages of development or not available. 13 Further, the large-scale data required to evaluate system performance cannot easily be aggregated 14 across charging networks. Given the rise in real-time streaming data in transportation and mo-15 16 bility apps, much of the useful intelligence about charging infrastructure performance lies highly unstructured. Consumer reviews, for example, can be collected instantly from thousands of users, 17 and manually processing or analyzing this unstructured data to collect useful information has not 18

19 been possible. We argue that real-time, streaming data will be increasingly important for research 20 evaluation of sustainable infrastructure. For example, since the release of a popular charging sta-21 tion locator app, there have been over 1.5 million user reviews of charging stations lying dormant as 22 text (7). Given this volume of data, even at an expert processing rate of 120 reviews per hour, it will 23 be prohibitively costly for humans to classify this unstructured text data for research evaluation.

We have shown in prior research that consumer sentiment can be automatically processed 24 25 with high accuracy through computational aid (8). From this analysis, the evidence suggests that there is a significant amount of negative consumer sentiment related to the charging experience. 26 While we demonstrated state-of-the-art performance from neural network-based models in this 27 domain, learning about the sources of negative consumer experience, which is needed to conduct 28 policy analysis, remains a challenging task due to the complexity of natural language processing 29 (NLP). In this paper, we therefore introduce a computational solution to analyze the content of real-30 time text data as tailored to the domain of EVs and consumer behavior. Because it is known that 31 consumer reviews may be subject to self-selection and other observational biases, we use machine 32 learning as a pre-processing tool to conduct econometric analyses for statistical adjustment. We use 33 this approach to automatically learn about large-scale barriers to EV infrastructure use nationally. 34 35 Contrary to the public discourse, we find that range anxiety is not a major concern among existing EV drivers. Instead, our results suggest that the major issues facing EV drivers relate to 36

station functionality and availability—an insight that we make possible through large-scale data
 integration.

39 MACHINE LEARNING AND TEXT CLASSIFICATION

40 With the increasing popularity of social data from digital platforms, user-generated short texts have

41 become an important data source for NLP. In the transportation domain regarding EV adoption,

- 42 there are as yet few research studies that translate the unstructured data from user generated texts 43 into actionable intelligence. One exception is a recent paper by Kuhl et al. (2019) in which
- 44 the authors manually coded Twitter data and found that contrary to what has been the focus in the

literature, charging infrastructure was the most discussed topic (9). There is no definitive study that 1 does large-scale analysis of EV behavior for policy analysis. Further, the implementation of recent 2 3 advances, such as deep neural networks (DNNs), which have revolutionized the field of natural language processing (10) have not been implemented for transportation policy. There are two 4 prevailing types of DNNs. The first are convolutional neural networks (CNNs). CNNs extract the 5 most meaningful information from text data by decomposing the hierarchical structure of sentences 6 or phrases (11). The other prevailing DNN architecture is that of recurrent neural networks (RNNs) 7 (12). RNNs have the added benefit of flexibility in analyzing a sequence of text. For instance, a 8 recent study by Ma et al., (2019) collected Chinese consumers' online comments about electric 9 10 vehicles and processed the data with a variant of an RNN architecture, known as the long shortterm memory (LSTM) in order to review EV purchase preferences such as retail prices, and EV 11 makes and models (13). However, the authors do not report their classifier performance measures 12 in their text mining analysis (e.g., accuracy and F1 score), so it is not possible to meaningfully 13 access the relative merits of DNNs in this domain. 14

Choosing between CNN or RNN architectures for any type of data is currently an ongo-15 16 ing debate in the literature (10). In our prior work, we initially demonstrated that a CNN model 17 produced state-of-the-art accuracy with good balance measures for sentiment classification tasks, which weakly dominated an RNN model (14). From the consumer analysis of charging station 18 19 sentiment, we were also able to demonstrate that there were differences in station-level sentiment when looking at the geographical regions (e.g. urban, rural, etc.) even after controlling for ob-20 servable station characteristics. However, sentiment analysis, although informative about quality 21 perceptions, does not give us a window into the specific mechanisms or sources of the negative 22 23 sentiment. As such, in this contribution, we are interested in applying innovations in NLP to introduce a typology of charging behavior and to provide insights on the use of computational tools for 24 25 the automatic detection and discovery of barriers to infrastructure management.

26 DATA AND METHODS

We have a nationally representative sample of unstructured consumer reviews at 12,720 US charg-27 ing station locations as provided by a popular EV charge station locator app. The text data consists 28 of 127,257 reviews written in English from 29,532 registered and unregistered EV drivers during 29 the period from 2011 to 2015. The sample represents charging stations from the entire U.S. market 30 during the period of study. This includes data aggregated from 10 major EV charging networks 31 in the US. In the sample, we also geocoded point of interest (POI) location information using 32 Google places API for categories such as Dealerships, Government, Healthcare, Hotel/Lodging, 33 Other, Park, Parking Garage/Lot, Residential, Restaurant, School/University, Shopping Center, 34 Store/Retail, and Workplace. For more information, ref. (8). 35

36 EV Charging Infrastructure Consumer Reviews

The charging station reviews can be considered social interactions within the community of EV drivers. After analyzing the contents of over 8,000 reviews, two research assistants were able to

 $\frac{1}{100}$ $\frac{1}$

- 39 identify the main categories discussed by users and were able to determine which issues are most 40 prominent regarding the charging experiences of this community. In preliminary experiments,
- 40 prominent regarding the enarging experiences of this community. In preminary experiments, 41 we investigated several unsupervised topic modeling techniques that did not provide theoretically
- 42 meaningful clusters. Therefore, we took the approach of hard coding labels based on human intel-
- 43 ligence. We introduce 9 main categories and 34 subcategories that make up a typology of charging



FIGURE 1: Map of Charging Stations in North America

behavior that allows for easier identification and eradication of inefficiencies of the charging pro-1 cess. The typology is provided in Table 1. Functionality refers to comments describing whether 2 particular features or services are working properly at a charging station. Range Anxiety refers to 3 comments regarding EV drivers' fear of running out of fuel mid-trip and to comments concerning 4 tactics to avoid running out of fuel. Availability refers to comments concerning whether charg-5 ing stations are available for use at a given station. Cost refers to comments about the amount of 6 money required to park and/or charge at particular locations. User Interaction refers to comments 7 in which users are directly interacting with other EV drivers in the community. Location refers 8 to comments about various features or amenities specific to a charging station location. The Ser-9 vice Time category refers to comments reporting charging rates (e.g. 10 miles of range per hour 10 charged) experienced in a charging session. The Dealerships category refers to comments con-11 cerning specific dealerships and user's associated charging experiences. The Other category refers 12 to comments that do not fall into the previous eight categories. From our sample of human labeled 13 reviews, the Other category occurs 6.0% of the time. For the full frequency counts by label, see 14 Table 2. 15

16 Approach to Curating the Training Data

17 Classification techniques employed by many scholars often assume that observational data is a

18 random sample from a given distribution that is believed to be representative of the population.

19 However, well-known biases in learning and evaluating classifiers can include researcher bias,

20 sample selection bias, and other statistical sampling issues (15-17).

For this reason, we actively curated a population of human annotators that were pre-

screened to be representative of the US general population (age 18+). With the support of aQualtrics panel, a sample of 1,000 participants were recruited with nationally representative demo-

24 graphic characteristics such as age, income and education level, sex and ethnicity. This allowed us

to explore techniques to mitigate, although not completely eliminate, potential individual bias by

Category	Subcategory
Functionality	General Functionality, Charger, Screen, Power Level, Connector Type,
	Card Reader, Connection, Time, Error Message, Station, Mobile Application,
	Customer Service
Range Anxiety	Trip, Range, Location Accessibility
Availability	Number of Stations Available, ICE, General Congestion
Cost	Parking, Charging, Payment
User Interactions	Charger Etiquette, Anticipated Time Available, User Tips
Location	General Location, Directions, Staff, Amenities, Points of Interest
	User Activity, Signage
Service Time	Charging Rate
Dealership	Dealership Charging Experience, Competing Brand Quality,
	Relationship with Dealers
Other	General Experiences

TABLE 1: EV Mobile App Typology of User Reviews

eliciting the wisdom of crowds and reducing the potential impact of conflicting biases. 1

2 We deployed an online survey questionnaire from Nov 27 to Dec 11, 2018 to build a training dataset for supervised machine classification. Each participant was given a total of 20 charging 3 station reviews and was tasked with labeling the sentiment of the review (positive or negative), 4 as well as selecting categories and sub-categories that applied to the review. Of the 20 charging 5 station reviews labeled by each participant, 5 of them were sample of 830 reviews designed for re-6 liability checks, while the remaining 15 reviews were randomly selected from the superpopulation. 7 The sample 830 reviews for the reliability checks were randomly drawn from a set that had been 8 previously labeled by research assistants. The holdout samples were designed to be distributed to 9 10 at least three of the 1,000 recruited participants on average for the purpose of calculating inter-rater agreement. All reviews that include potentially sensitive information such as cell phone numbers 11 or email addresses were redacted. Because subcategories of many categories such as Function-12 ality had domain-specific terminology that could be unfamiliar to a general population, we also 13 provided a diagram of a charging station with information along with tips to help guide the human 14

classification. A view of the online survey interface is shown in Figure 2. After the labeling tasks 15

were completed, we also asked participants for voluntary demographic information. No personally 16

identifiable information was collected or shared (IRB protocol No. H18250). 17

Inter-Rater Reliability 18

Because we had multiple annotators, we used Fleiss' Kappa (κ) as a measure of agreement be-19 tween raters. Fleiss' Kappa was selected because it offers the benefit of a single metric to assess 20

21 agreement between *n*-raters (18). On average, we had an average of 3 raters per review, ranging

22 from 1 to 7 in the experiment. The Fleiss' Kappa, κ is calculated as below:

$$\kappa = \frac{\bar{P} - \bar{P}_e}{1 - \bar{P}_e},\tag{1}$$

23 where \bar{P} is the average number of agreement on all category assignments between rater pairs for 24

the questions, and \bar{P}_e is the average proportion of assignment to the categories. For example,



FIGURE 2: Online Survey Interface

- 1 assume there are 3 raters labeling 3 questions with 3 categories each. Suppose that the 3 raters
- 2 agree on 2 questions with distinct categories and completely disagree on 1 question, choosing non-
- 3 overlapping categories. In this case, $\vec{P} = \frac{2}{3}$, and $\vec{P}_e = \frac{11}{27}$, resulting in a κ score of 0.44. As κ is 4 bounded between -1 and 1, when κ is less than 0, agreement between raters is occurring below
- 4 bounded between -1 and 1, when κ is less than 0, agreement between raters is occurring below 5 what would be expected at random while a κ above 0 means that agreement between raters is
- 6 occurring above what would be expected at random while a k above 0 means that agreement between raters is
- 7 would interpret the κ as moderate agreement. For more information, see ref. (18). We provide an
- analysis of inter-rater reliability in the Results and Discussion section.

9 CNN Implementation

- 10 In this study, we expand on CNN implementation protocols introduced in ref. (14), which is part
- 11 of a larger literature of consumer analysis using DNNs (20–22). Given that many labels can occur
- 12 relatively infrequently leading to unbalanced class data for training and testing, a simple majority
- 13 classifier that assigns the majority class to all test data is provided as a baseline for performance.
- 14 The performance of the CNN classifier is also compared to other commonly used baseline classifi-
- 15 cation models, e.g. the bag-of-N-grams based logistic regression (LR) and support vector machines
- 16 (SVM).

17 Convolutional Neural Network for Text Classification

- 18 Here we provide a high level overview of a single layer CNN model for binary topic classification.
- 19 Figure 3 shows a generalized CNN architecture, which is comprised of three main parts: the vector
- 20 representation of words, the convolutional layer, and the fully connected layer that yields a binary
- 21 category prediction as the output. First, all words in a consumer review are changed to tokens,
- 22 each having a certain dimensional vector composed of numeric values. These values are called
- 23 word embeddings, which represent similarity between vocabularies quantitatively; similar words

have similar values across their word vectors and vice versa. In this study, we used pre-trained, 1 publicly available word vectors, namely the *word2vec* (23), which is a predictive model trained 2 3 on 100 billion words and phrases in Google News. We also tried other open-sourced, pre-trained word vectors, GloVe (24), a count-based model which gave comparable results. The use of pre-4 trained word vectors in text classification task is a widely used method to improve performance in 5 the absence of a large corpus of domain-specific training data (20, 25). In the convolutional layer, 6 filters scan through the word vectors creating what are referred to as "feature maps" using activa-7 tion functions that are representative of the word vector of the chosen filter size. This process is 8 done empirically by multiple filters, and the results are combined and extracted in a process called 9 "pooling." Extracted information by pooling process is gathered as one feature vector, and it cre-10 11 ates prediction through the fully connected layer. A CNN model can have various hyperparameters for the process of filtering, convolution, and pooling processes. For more information about details 12 of the algorithm, see refs, (11, 20). 13



FIGURE 3: Architecture of the Convolutional Neural Network Model

- 14 Training and Testing Data
- 15 We created nine binary classifiers for each of the main categories: Functionality, Range Anxiety,
- 16 Availability, Cost, User Interactions, Location, Service Time, Dealership and Other. For the train-
- 17 ing data, we used the 20,000 reviews labeled by the 1,000 participants. The models are validated
- 18 by a set 5,229 true labels provided by research assistants as human experts. The counts of labels
- 19 provided by participants are summarized in the Table 2. As the categories labeled in the training
- 20 set are largely imbalanced, it is important to evaluate balance measures in classifier performance
- 21 and to verify the level of learning compared with the simple majority class model. Based on the
- 22 most common labels, the most important topics to consumers are *Functionality* and *Availability*.

	Labeled	Not Labeled	Percent Labeled ^{\dagger} (%)
Functionality	7,370	12,630	36.9
Range Anxiety	1,809	18,191	9.0
Availability	4,968	15,032	24.8
Cost	2,006	17,994	10.0
User Interactions	3,174	16,826	15.9
Location	3,046	16,954	15.2
Service Time	1,805	18,195	9.0
Dealership	1,075	18,925	5.4
Other	1,209	18,791	6.0

TABLE 2: Counts of Labeled Reviews per Category in Training Data

[†] Reviews can have multiple labels therefore Percent Labeled does not sum to 100%.

1 Hyperparameter Optimization

2 We perform basic hyperparameter tuning, beginning with suggested values introduced in previous

3 studies (14, 20, 26). These include filter region sizes of [3, 4, 5]; 100 filters; use of rectified linear

4 unit (ReLU) activation function; learning rate of 0.001; dropout rate of 0.3 for regularization, and

5 no l_2 norm constraint. We also followed guidelines provided by (26) that each dataset has its own

6 optimal filter region size and numbers, therefore perfromed basic grid search to find hyperparame-

7 ters that could yield higher performance in our specific dataset. For the Functionality category, we

8 found that filter region sizes of [12, 12, 12], 400 number of filters, learning rate 0.0001, 0.6 dropout

9 rate, 3 epochs with 128 batch size resulted in the highest performance within our search. However,

10 with other categories, the performance dropped as changes occurred to our initial hyperparameters.

11 Outcomes of Interest

12 We are interested in evaluating the factors that predict the performance of outcomes across stations.

13 Given our objective of detecting behavioral failures, we focus our analysis on *Functionality* and

14 Availability, which are the 2 most frequently observed labels in the training data. We created an

15 index that measures the probability that a label is likely to be chosen. For a given station review i, 16 at leastion group a in ware the Label Score is defined as follows:

Label Score_{*m,i,g,year*} =
$$\frac{\text{Count of label reviews}_{m,i,g,year}}{\text{Total count of reviews}_{i,g,year}}$$
 (2)

17

18 where m is the label of interest. A score near 0 would indicate that the specific label has a low

19 incidence at that location group in that year, while a high score near 1 would indicate a high

20 incidence at the location group in that year were assigned that label by the classifier.

21 Econometric Estimation using Fractional Response Models

22 The most commonly used implementation for models with fractional dependent variables only

23 requires some specification of the correct functional form of the fractional dependent variable

24 (27). In our case, the variables of interest are fractional vectors bounded between 0 and 1 with

25 probability masses accumulating at the boundaries of the interval (See Figure 4). It would be



FIGURE 4: Histograms of Label Scores

- 1 inappropriate to use an OLS or log-odds estimator, for reasons we have described elsewhere (14).
- 2 In the general form for the fractional response model, the interest is on the conditional expectation
- 3 of the fractional response variable $y_{i,t}$ on the vector of explanatory variables $\mathbf{x}_{i,t}$ such that,

$$E(y_{i,t}|x_{i,t}) = G(\mathbf{x}_{i,t}\boldsymbol{\theta}), \ i = 1,\dots,N,$$
(3)

4

5 where $G(\cdot)$ is a non-linear transform function where the cdf satsifies $0 \le G(\cdot) \le 1$, the fractional 6 dependent variable is defined only on $0 \le y_{i,t} \le 1$, and θ is a parameter vector of interest. Estimates 7 of the effects are directly computed using the Permeulli log likelihood function given by (28, 20)

7 of the effects are directly computed using the Bernoulli log-likelihood function given by (28, 29),

$$LL_{i,t}(\boldsymbol{\theta}) \equiv y_{i,t} log[G(\mathbf{x}_{i,t}\boldsymbol{\theta})] + (1 - y_{i,t}) log[1 - G(\mathbf{x}_{i,t}\boldsymbol{\theta})].$$

$$(4)$$

9 Given the presence of boundary observations of 0 in the fractional dependent variables of our 10 dataset, the pooled Bernoulli quasi-maximum likelihood estimator (QLME) of θ can be computed 11 as,

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{N} LL_{i,t}(\theta).$$
(5)

13

For our analysis, the main functional specification regresses the outcomes of interest on a vector of dependent variables, including geographical area dummies, station characteristics, POI dummies, Negativity Score, and related interactions. For the geographical area dummy variables we utilize the Census definitions for urban centers, urban clusters, and rural areas. From the Census definitions, urban centers are urbanized areas with more than 50,000 people, urban clusters are

- 1 urbanized ares with populations between 2,500 and 50,000 people, and rural areas are any such
- 2 areas outside of urbanized ones (30). The station characteristics include measures of the number
- 3 of connectors at a station, the number of networks at a station, and a proprietary station quality
- 4 rating that ranges between 1 and 5. The general functional specification is as follows:

Label Score_{*m,i,g,year*} =
$$\alpha_{i,year}$$
 + Geographical Area Dummies_{*g*} + Station Characteristics_{*i,g,year*}
+ Negativity Score_{*i,g,year*} + Interaction Effects_{*i,g,year*} + POI Dummies_{*g*} (6)

5 RESULTS AND DISCUSSION

6 Descriptive Analysis of the Training Data

- 7 We performed a basic analysis to learn about the characteristics of the collected data. We provide
- 8 a frequency distribution of the occurrences of all hard-coded labels in Figure 5 for main categories
- 9 and in Figure 6 the most frequently occurring subcategories for our top category. The most selected
- 10 main categories were Functionality and Availability, with counts of 7,370 and 4,968 respectively.
- 11 We also learned that *Range Anxiety* is not a dominant theme among EV drivers' discussions (only
- 12 9% of the training data). We also break out the Functionality subcategories which were most dis-
- 13 cussed. For example, many consumers report issues with the proper functioning of the charger
- 14 itself, station features, customer service, wireless connection, power issues and list of other prob-
- 15 lems, see Figure 6.



FIGURE 5: Frequency Count of Training Data per Main Category

Using the holdout sample labels collected from the training data, the Fleiss' κ was measured to be $\kappa = 0.31$, which is considered "fair agreement" in prior literature (19, 31). Although we note that there are relatively few existing studies that experimentally curate the training data, or report their reliability measures, we found some recent studies that begin to report reliability measure from human annotators using social data from posts in online forms (32). Further details on the effects of inter-rater reliability on classifier performance are discussed in the next section.



FIGURE 6: Frequency Counts per Subcategory of Functionality.

1 Classification Results

2 We are interested to learn whether a neural net-based model can accomplish the task of detecting

3 a performance category from text with good accuracy compared to humans, while still achieving

4 good balanced measures, using the F1 score as an indicator. For the analysis of classifiers we

5 report model metrics for our top 3 categories which are *Functionality*, *Cost*, and *Availability*. Table

6 3 shows the comparison of CNN versus the other baseline models.

	Accuracy (%)			F1 Score				
	CNN	SVM	LR	Majority	CNN	SVM	LR	Training Data
Functionality	78.1	73.3	74.5	54.7	0.72	0.66	0.66	Balanced
Cost	94.1	93.3	92.9	90.6	0.59	0.57	0.44	Highly
Availability	85.2	83.2	0.85	88.7	0.50	0.37	0.46	Imbalanced Moderately Imbalanced

TABLE 3: Classification Results for Top Categories

Despite the fact that we have imbalanced training data, we see that the CNN model outper-7 formed the three baseline models for most categories, and is significantly better at detecting true 8 9 positives and true negatives as shown by the F1 score. For example, one review states "Right not 10 working. Terrible". Our human annotators identified this review as a Functionality category. For the machine classification, LR and SVM failed to label Functionality, while CNN learned from 11 semantic context to properly identify the label. In the opposite context of misclassification, CNN 12 also performed better. For example, a consumer writes "Stephen, you can use the Greenlots app 13 if you don't have a fob.". LR and SVM both predicted this to be Functionality label, while both 14 human annotators and the CNN classifier correctly identified this text as not relating to Function-15 ality. These examples illustrate that the CNN model has strong potential to automatically learn 16 17 about the major issues in this domain. The results we report in Table 3 represent state-of-the-art 1 classification performance known for this domain.

2 One exception to the strong performance of CNN is the Availability category, where the 3 accuracy of CNN was slightly lower than the simple majority class algorithm. This may be due to the difference in labeling rules between participants and trained annotators, where in the training 4 data, Availability was labeled for 24.8% of the training reviews, while the trained annotators were 5 more strict and labeled Availability for only 11.3% of the testing reviews. This suggests that in 6 future work, we could provide more clear instructions to general population of annotators about the 7 definitions of the labels. We also found that the Cost category had better F1 measure as compared 8 9 ton Availability with higher imbalance in the training data. 10 However this strong performance is not without limitations. It did not do well on the labels

such as Location, Service Time, User Interactions, Dealership, Range Anxiety, and Others. Based 11 on our survey questionnaire, our results suggest that many of these domain specific terms can be 12 complex for the general US population. To investigate this further, we calculated the Fleiss' κ 13 of 0.31, which indicates that the inter-rater reliability score can be significantly enhanced. For 14 this reason, in future work, we suggest curating crowd sourced human labels for machine learning 15 using crowds of experts who might be more proficient in this domain. We leave that as future 16 work. Further, we have implemented one neural net-based model, but it is not the only possible 17 architecture. In addition to further tuning of the model we presented in this study, we also suggest 18 exploring other deep learning architectures, particularly the recurrent neural networks which could 19 have the benefit of learning from sequences of text data. 20

Having shown the good performance of our classier, we then use this model as a preprocessing tool to conduct econometric analyses in order to evaluate large-scale consumer issues in charging infrastructure. For more details on model performance, results are available upon request.

25 Fractional Response Model Results

We use the training data to classify the full population of 127,257 reviews in order to focus on our top two categories for the fractional response models (FRMs). Our main objectives are to uncover the main categories driving negative experience in the charging infrastructure in the United States. We evaluate geographical areas, both urban and non-urban, and how explanatory features are moderated by negative experience automatically classified using machine intelligence. The main results of the FRMs can be seen in Table 4.

32 Negative Charging Experiences with Functionality More Likely in Urban Center

33 In model (I) we estimate a basic specification that regresses the Functionality Score on observable station characteristics. We cluster our standard errors at the location group level. We find no 34 evidence that geographical area is a significant predictor of functionality labels; however, in models 35 (II) and (III) we investigate the sub-population of reviews with high Negativity Scores and find 36 that reviews urban centers are 46.6% more likely to be about Functionality. For example, one 37 user posted a review about Functionality in an urban center that state, "one of two chargers has 38 fault error and doesn't work both charge a fee:(", while the same user posted a review outside of 39 an urban center that stated, "120v outlet on lamp by entrance/sign very nice hotel excellent food 40 next door", which is not about Functionality. This result answers a previously open question about 41 the sources of negative consumer sentiment. By contrast, urban clusters and rural areas were not 42 significant predictors of *Functionality* labels net of all controls. 43

Ha, Marchetto, Burke and Asensio

1 The models used were robust to the inclusion of the quality rating as well as to various 2 clustering alternatives for standard errors. These additional results are available upon request.

3 Negative Charging Experiences with Availability More Likely in Urban Cluster

We also evaluated factors related to station availability. In model (IV), we estimate a basic spec-4 5 ification that regresses the Availability Score on all observable station characteristics. We do find evidence that urban centers have more reviews relating to Availability topics as compared to rural 6 areas. Our results show that the sources of negative consumer experiences related to Availability 7 are primarily in the urban clusters with reviews being 39.6% more likely. By contrast we do not 8 find statistically significant results in with negative consumer sentiment in urban center or rural 9 areas. This is interesting because one would expect Availability issues to be an urban phenomenon 10 related to congestion. However, our large-scale analysis points to smaller urban clusters could need 11 additional resources to broaden the availability of charging stations. 12

	Fur	nctionality S	Score	Availability Score			
	(I)	(II)	(III)	(IV)	(V)	(VI)	
Geographical Area							
Urban Center	-0.080	-0.273**	-0.280**	0.455***	0.457***	0.461***	
	(0.094)	(0.132)	(0.132)	(0.091)	(0.091)	(0.091)	
Urban Cluster	0.185	0.179	0.188	-0.243**	-0.399***	-0.426***	
	(0.153)	(0.147)	(0.148)	(0.122)	(0.147)	(0.149)	
Station Characteristics							
Number of Connectors	0.450***	0.451***	0.456***	-0.477***	-0.476***	-0.482***	
	(0.025)	(0.025)	(0.025)	(0.037)	(0.037)	(0.037)	
Number of Networks	-0.148*	-0.147*	-0.162*	0.179	0.180	0.177	
	(0.085)	(0.085)	(0.088)	(0.154)	(0.154)	(0.156)	
Quality Rating			-0.059***			0.075***	
			(0.017)			(0.017)	
Negativity Score	1.744***	1.339***	1.246***	0.813***	0.798***	0.901***	
	(0.054)	(0.220)	(0.216)	(0.055)	(0.056)	(0.055)	
Urban Center x Negativity Score		0.466**	0.474**				
		(0.227)	(0.224)				
Urban Cluster x Negativity Score					0.396**	0.433**	
					(0.201)	(0.206)	
Point of Interest Control Dummies	Yes	Yes	Yes	Yes	Yes	Yes	
Clustered SE	Yes	Yes	Yes	Yes	Yes	Yes	
Number of Observations	127,257	127,257	127,257	127,257	127,257	127,257	
R2	0.154	0.155	0.158	0.053	0.053	0.056	
Note:				* <i>p</i> <0.1;	** <i>p</i> <0.05;	***p<0.01	

TABLE 4: Fractional Response Model Results

Among the observable station characteristics, it turns out the the number of connectors is a significant predictor of both functionality and availability topics. In the case of *Availability*, more connectors predicts less reviews with availability labels. With *Functionality*, more connectors predicts more reviews with functionality labels. This is intuitive as more connectors should help availability issues while also providing more chances for functional issues to occur at a charge station. Overall, sources of negative consumer sentiment appear to be an urban phenomenon with urban centers being related to negative sentiment in reviews about functionality and urban clusters 1 being related to negative sentiment in reviews about availability.

2 POLICY IMPLICATIONS

In this study, we have been able to demonstrate that advances in computational algorithms can 3 be deployed at relatively low cost with promising performance measures. We also demonstrate 4 that through large-scale data aggregation, it may be possible to build a framework that could cap-5 ture consumer intelligence about the functioning of the infrastructure, in near-real-time. Such 6 capabilities could revolutionize how we manage, evaluate, and invest in charging infrastructure for 7 electrified transportation. Based on our results, we have three main policy recommendations. First, 8 a necessary criteria for building frameworks for real-time analysis is data and information sharing. 9 We suggest the expansion of policies that can allow for greater real-time data sharing regionally 10 and between jurisdictions. Second, given the discovery of negative consumer experiences as an 11 urban phenomena, we suggest strategies for local and regional government to push greater stan-12 dards and investment to help ensure that the quality and reliability of the charging experience is 13 core to policies for EV growth. Third, the discussion of range anxiety as consumer barrier appears 14 to be overstated, whereas station functionality and charging availability at the point of use may be 15 the more critical limitation. Real-time streaming data is already changing the nature of mobility 16 decisions for consumers. With a trained model present, thousands of reviews can be processed 17 and analyzed in matter of minutes, if not faster. This should yield significant cost reductions for 18 infrastructure performance evaluation. 19

20 ACKNOWLEDGEMENTS

- 21 We thank the generous support of the National Science Foundation Award No. 1931980, the An-
- 22 thony and Jeanne Pritzker Family Foundation, Microsoft Azure for Research (Award CRM:0518988),
- 23 the Sustainable LA Grand Challenge, and the Ivan Allen College Dean's SGR-C Award. For valu-
- 24 able research assistance, we thank Cade Lawson and Soobin Oh. We thank Suzie Lee and Han-
- 25 nah Weirich at Qualtrics for participant support. For valuable comments and feedback, we thank
- 26 Richard Fujimoto, Emily Grubert, Haesun Park and Iris Tien. This research was supported in part
- 27 through research cyber-infrastructure resources and services provided by the Partnership for an Ad-
- 28 vanced Computing Environment (PACE) at the Georgia Institute of Technology, Atlanta, Georgia,
- 29 USA.

1 **REFERENCES**

- EPA, Inventory of U.S. Greenhouse Gas Emissions and Sinks: 1990-2016, 2018, document
 No. 430-R-18-003.
- Council, N. R. et al., *Hidden costs of energy: unpriced consequences of energy production and use.* National Academies Press, 2010.
- Bepartment of Energy, *Electric Vehicles: Tax Credits and Other Incentives Database*,
 2019, access date: 07/31/2019, https://www.energy.gov/eere/electricvehicles/
 electric-vehicles-tax-credits-and-other-incentives.
- 9 4. Carley, S., R. M. Krause, B. W. Lane, and J. D. Graham, Intent to purchase a plug-in
 10 electric vehicle: A survey of early impressions in large US cites. *Transportation Research*11 *Part D: Transport and Environment*, Vol. 18, 2013, pp. 39–45.
- Sheldon, T. L., J. DeShazo, and R. T. Carson, Electric and plug-in hybrid vehicle demand:
 lessons for an emerging market. *Economic Inquiry*, Vol. 55, No. 2, 2017, pp. 695–713.
- Navigant, Market Data: EV Market Forecasts: Global Forecasts for Light Duty Plug-In
 Hybrid and Battery EV Sales and Populations: 2017-2016, 2017.
- Recargo, *PlugShare Key Features and Benefits*, 2019, access date: 07/31/2019, https:
 //recargo.com/plugshare.html.
- Alvarez, K., A. Dror, E. Wenzel, and O. I. Asensio, Evaluating Electric Vehicle User
 Mobility Data using Neural Network based Language Models. In *Proceedings of 98th Annual Meeting of the Transportation Research Board, ADC80 Standing Committee on Alternative Transportation Fuels and Technologies, Washington, D.C.*, 2019.
- Kuhl, N., M. Goutier, A. Ensslen, and P. Jochem, Literature vs. Twitter: Empirical insights
 on customer needs in e-mobility. *Journal of Cleaner Production*, Vol. 213, 2019, pp. 508–
 520.
- Yin, W., K. Kann, M. Yu, and H. Schütze, Comparative study of CNN and RNN for natural
 language processing. *arXiv preprint arXiv:1702.01923*, 2017.
- LeCun, Y., L. Bottou, Y. Bengio, P. Haffner, et al., Gradient-based learning applied to
 document recognition. *Proceedings of the IEEE*, Vol. 86, No. 11, 1998, pp. 2278–2324.
- 29 12. Elman, J. L., Finding structure in time. *Cognitive science*, Vol. 14, No. 2, 1990, pp. 179–
 30 211.
- Ma, S. C., Y. Fan, J. F. Guo, J. H. Xu, and J. Zhu, Analysing online behaviour to determine
 Chinese consumers' preferences for electric vehicles. *Journal of Cleaner Production*, Vol.
 229, 2019, pp. 244–255.
- Asensio, O. I., K. Alvarez, A. Dror, E. Wenzel, C. Hollauer, and S. Ha, Evaluating popular
 sentiment of electric vehicle owners in the United States with real-time data from mobile
 platforms. *Working paper*, 2019.
- Shepperd, M., D. Bowes, and T. Hall, Researcher bias: The use of machine learning in
 software defect prediction. *IEEE Transactions on Software Engineering*, Vol. 40, No. 6,
 2014, pp. 603–616.
- 40 16. Zadrozny, B., Learning and evaluating classifiers under sample selection bias. In *Proceed-*41 *ings of the twenty-first international conference on Machine learning*, ACM, 2004, p. 114.
- 42 17. Dwork, C., V. Feldman, M. Hardt, T. Pitassi, O. Reingold, and A. L. Roth, Preserving 43 statistical validity in adaptive data analysis. In *Proceedings of the forty-seventh annual*
- 44 ACM symposium on Theory of computing, ACM, 2015, pp. 117–126.

- Fleiss, J. L., Measuring nominal scale agreement among many raters. *Psychological bulletin*, Vol. 76, No. 5, 1971, p. 378.
- 3 19. Landis, J. R. and G. G. Koch, The Measurement of Observer Agreement for Categorical
 4 Data. *Biometrics*, Vol. 33, No. 1, 1977, pp. 159–174.
- 5 20. Kim, Y., Convolutional Neural Networks for Sentence Classification. *arXiv:1408.5882*6 [*cs*], 2014, arXiv: 1408.5882.
- Kalchbrenner, N., E. Grefenstette, and P. Blunsom, A convolutional neural network for
 modelling sentences. *arXiv preprint arXiv:1404.2188*, 2014.
- 9 22. Johnson, R. and T. Zhang, Effective Use of Word Order for Text Categorization with Convolutional Neural Networks. *arXiv:1412.1058 [cs, stat]*, 2014, arXiv: 1412.1058.
- Mikolov, T., I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, Distributed Representations of Words and Phrases and their Compositionality. In *Advances in Neural Information Processing Systems 26* (C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q.
 Weinberger, eds.), Curran Associates, Inc., 2013, pp. 3111–3119.
- Pennington, J., R. Socher, and C. Manning, Glove: Global Vectors for Word Representa tion. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Association for Computational Linguistics, Doha, Qatar, 2014, pp.
 1532–1543.
- Iyyer, M., J. Boyd-Graber, L. Claudino, R. Socher, and H. Daumé III, A Neural Network
 for Factoid Question Answering over Paragraphs. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Association for Computational Linguistics, Doha, Qatar, 2014, pp. 633–644.
- Zhang, Y. and B. Wallace, A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, Asian Federation of Natural Language Processing, Taipei, Taiwan, 2017, pp. 253–263.
- 27 27. Ramalho, E. A., J. J. Ramalho, and J. M. Murteira, Alternative estimating and testing empirical strategies for fractional regression models. *Journal of Economic Surveys*, Vol. 25, No. 1, 2011, pp. 19–68.
- Papke, L. E. and J. M. Wooldridge, Econometric methods for fractional response variables
 with an application to 401 (k) plan participation rates. *Journal of applied econometrics*,
 Vol. 11, No. 6, 1996, pp. 619–632.
- Papke, L. E. and J. M. Wooldridge, Panel data methods for fractional response variables
 with an application to test pass rates. *Journal of Econometrics*, Vol. 145, No. 1-2, 2008,
 pp. 121–133.
- 36 30. US Census, 2010 Urban Area FAQS, 2010.
- 37 31. Gwet, K. L., Handbook of inter-rater reliability: the definitive guide to measuring the
 as extent of agreement among raters; [a handbook for researchers, practitioners, teachers &
 students]. Advanced Analytics, Gaithersburg, MD, 3rd ed., 2012.
- 40 32. Bobicev, V. and M. Sokolova, Inter-Annotator Agreement in Sentiment Analysis: Machine
 41 Learning Perspective. In *RANLP*, 2017, pp. 97–102.