## Using Public Data to Predict Demand for Mobile Health Clinics

Haipeng Chen, <sup>12</sup> Susobhan Ghosh,<sup>12</sup> Gregory Fan,<sup>34</sup> Nikhil Behari,<sup>12</sup> Arpita Biswas, <sup>12</sup> Mollie Williams,<sup>34</sup> Nancy E. Oriol,<sup>34</sup> Milind Tambe<sup>12</sup>

 <sup>1</sup> Center for Research on Computation and Society
 <sup>2</sup> John A. Paulson School of Engineering and Applied Sciences, Harvard University
 <sup>3</sup> The Family Van <sup>4</sup> Harvard Medical School Contact: hpchen@seas.harvard.edu

#### Abstract

Improving health equity is an urgent task for our society. The advent of mobile clinics plays an important role in enhancing health equity, as they can provide easier access to preventive healthcare for patients from marginalized populations. For effective functioning of mobile clinics, accurate prediction of demand (expected number of individuals visiting mobile clinic) is the key to their daily operations and staff/resource allocation. Despite its importance, there have been very limited studies on predicting demand of mobile clinics. To the best of our knowledge, we are the first to explore this area, using AI-based techniques. A crucial challenge in this task is that there are no known existing data sources from which we can extract useful information to account for the exogenous factors that may affect the demand, while considering protection of client privacy. We propose a novel methodology that completely uses public data sources to extract the features, with several new components that are designed to improve the prediction. Empirical evaluation on a real-world dataset from the mobile clinic The Family Van shows that, by leveraging publicly available data (which introduces no extra monetary cost to the mobile clinics), our AI-based method achieves 26.4% - 51.8% lower Root Mean Squared Error (RMSE) than the historical average-based estimation (which is presently employed by mobile clinics like The Family Van). Our algorithm makes it possible for mobile clinics to plan proactively, rather than reactively, as what has been doing.

### **1** Introduction

The disproportionate impact of the COVID-19 pandemic on marginalized populations has exemplified long-standing health inequities in the United States (Mackey et al. 2021; Zimmerman and Anderson 2019; Odlum et al. 2020). These health disparities are often the result of marginalized populations facing increased barriers to healthcare access, including fear or mistrust of the medical system, and prohibitive travel times (Syed, Gerber, and Sharp 2013; Bolen et al. 2016). In response to both the pandemic and the health disparities, healthcare providers and public health organizations have begun implementing novel solutions, including the increased use of telecare and social media (Vance, Howe, and Dellavalle 2009; Woolliscroft 2020). One such solution, the mobile health clinic, is a large bus or van that has been

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

converted to provide medical care at a location physically and socially closer to at-risk communities (see Figure 1 as an example of the studied mobile clinic The Family Van).<sup>1</sup> These clinics attempt to address multiple barriers to healthcare access by reducing patient travel times and focusing on providing culturally-competent care (Stephanie et al. 2017; de Peralta et al. 2019).



Figure 1: Clients visiting The Family Van in Boston. Photo source: The Family Van's website.

A main benefit of using mobile clinics is their ability to relocate to areas of high demand. Therefore, it is critical for mobile clinic administrators to understand what key factors affect client demand and what future demand might look like. Predictive demand models could then be used to optimize van scheduling as well as the allocation of staff and healthcare resources. Despite its importance, only a limited number of studies discuss forecasting client demand for mobile clinics. Existing works use either static features from census data (Majeed et al. 2021) which do not reflect demand *dynamics*, or patient survey data (Reed et al. 2019) which may violate patient privacy requirements, suffer from selection bias, or may not be available.

However, developing such demand prediction models involves several unique challenges: i) It is not clear which factors affect demand. ii) Unlike hospitals, data from mobile clinics are not as available. There are no known existing data sources from which we can extract useful information for exogenous factors that affect clinic demand while also protecting patient privacy. iii) For the factors that are identified as good indicators of demand, we have no prior knowledge

<sup>&</sup>lt;sup>1</sup>http://www.familyvan.org/

of their ground truth values on the forecast date. These challenges prevent us from directly applying off-the-shelf Machine Learning (ML) models.

We propose a novel prediction framework that addresses the above challenges, with the following innovations. (i) We find features from public data sources that we hypothesize to be good indicators of demand. We determine which factors should be included in the prediction model by performing a correlation analysis of each factor with respect to mobile clinic demand. Very interestingly and surprisingly, we find that factors such as ferry and shared bike usage, which are proxies of foot traffic, are highly correlated with client demand. (ii) Though we can gain insights from the correlation analysis, a major obstacle to using these insights to extract features for the prediction model is that the values of these features on the forecast date are not known a-priori. Therefore, we propose to first make intermediate predictions for the future values of these variables, and then use the predicted variables as features. (iii) We observe that the demand patterns of repeat/non-repeat clients are distinct. To further improve prediction, we separate out the demand predictions for each type of client, and then combine the predictions. Finally, we integrate these components with different ML models to make the final prediction. Whereas we do not claim novelty for these ML models, our method is an innovative combination of the known ML models with the above mentioned components.

Our main contributions are as follows: (i) Impact. We are the first to develop ML tools to dynamically predict the weekly/daily client demand of mobile health clinics. Except for the demand data (which is accessible to parties in the field), our model uses publicly available data to extract features. Therefore, our model can be used by a vast range of mobile health clinics, and have a huge impact on the mobile health industry, where making demand predictions is essential to proactive planning of their daily operations and investment strategy. (ii) Technical novelty. We propose a novel prediction model which has several novel components that are designed to address the new challenges in the demand prediction task of interest. (iii) Effectiveness. We run experiments using real-world data from 4 locations of the mobile health clinic The Family Van in the Boston area. Results show that our new AI-based approach has a 26.4% - 51.8% lower RMSE than the traditional practice of The Family Van.

### 2 Related work

We provide a brief discussion of related work on health care demand forecasting. Majeed et al. (2021) focus on predicting demand for mobile clinics that provide free vaccination services in schools and regional areas that are at high risk of infection. They estimate the demand using non-temporal (static) census features and school-level data. In contrast, our work considers temporal features and makes real-time forecasting possible. Reed et al. (2019) strictly use patient survey data and focus only on the perception of clinic quality. The paper examines associations between multiple features, such as perceived quality of care, travel distance, and patient utilization of alternative health care clinics, using multiregression models on the survey data. Although the perception of alternate health care has an impact on the demand of mobile clinics, but relying on patient survey data may violate patients' privacy requirements. Moreover, such findings may suffer from selection bias. Qian et al. (2009) focus on the factors that influence the health care demand of public and private clinics in rural areas of Gansu province, China. However, they do not consider mobile health clinics or temporal features. Similarly, there are several papers that use multiple regression models for demand/utilization prediction of traditional hospitals and healthcare centers using hospitalization rates and ICU beds (Bhowmik and Eluru 2021), radiology volume records (Côté and Smith 2018), surgeries and admissions volume (Khaldi et al. 2017), and blood donation records (Drackley et al. 2012). However, none of these studies can be directly used or mildly adapted for the task of predicting the demand of mobile clinics. As far as we know, our study is the first ML-based method to predict dynamic weekly/daily patient demand of mobile health clinics.

### **3** Demand prediction for mobile clinics

We first formalize the demand prediction problem for our studied mobile clinic, The Family Van, and then describe the challenges in solving it. The Family Van is a non-profit mobile health clinic in the Boston area, designed to increase access to health care and improve the health of Boston's most under-served neighborhoods. It sends out medical vans to 4 major locations in the Boston area, one day per location, thus operating 4 days per week. Our goal is to predict the demand of the four locations for a future time (e.g., next week), given the historical demand data.

**Problem formulation** Abstracting from the above scenario, we are given a time series  $1, \ldots, t, t + 1, \ldots$ , where each t is a week. There are various locations  $i = 1 \ldots N$ , one for each day of a week t. The demand at time t and location i is denoted as a non-negative integer  $y_{i,t} \in \mathbb{Z}^+$ . Suppose we are at time t - 1, our goal is then to predict the mobile clinic's demand  $y_{i,t}$  for each of the N locations, at the next time period t. Essentially, we want to learn the following demand function  $f_{\theta}(\cdot)$ :

$$\hat{y}_{i,t} = f_{\theta}(y_{i,1}, \dots, y_{i,t-1}), \ \forall i = 1 \dots N.$$
 (1)

 $f_{\theta}(\cdot)$  can be interpreted as a certain ML model. In the next section, we will introduce our choice of the model.

**Challenges** This prediction task is challenging because of three aspects. First, it is unknown which exogenous factors affect mobile clinics' demand and to what extent. It remains a question whether to introduce exogenous factors to the ML model, or merely use a time-series model. Second, unlike regular hospitals, the mobile clinic data are not as accessible (partly due to the need to protect patient privacy). There is no known existing data source for the exogenous factors that may affect demand. Last, for the factors that are identified as good indicators of demand, we have no prior knowledge of their ground truth values of the forecast date.



Figure 2: Architecture of our proposed prediction method. The main innovations are 1) using completely public data sources to obtain features, especially the bike and ferry usage features which are surprisingly good indicators of demand dynamics, 2) the intermediate predictions component, and 3) separate predictions for repeat/non-repeat clients.

## 4 Methodology

We now introduce the techniques for addressing the challenges. We refer to Figure 2 for a high-level illustration of our solution architecture.

### 4.1 Location specific model

Before determining the right ML model for prediction, our first observation, as shown in Fig. 3, is that the demand curves of different locations are distinct from each other. For example, the demand scales of different locations are different – the demand of location 1 can reach 30+, whereas the demand of location 2 is mostly under 15. In addition, the demand trend of location 3 is generally decreasing, whereas such a trend is not obvious for the other locations. Because of this, though a uniform and generalized model that works for all locations is desirable from a technical perspective, we choose to use a location specific model instead. In this case, Eq.(1) can be re-written as:

$$\hat{y}_{i,t} = f_{\theta_i}(y_{i,1}, \dots, y_{i,t-1}), \ \forall i = 1 \dots N.$$
 (2)

In other words, we will learn one separate model  $f_{\theta_i}(\cdot)$  for each location  $i = 1 \dots N$ .

## 4.2 Time-series model

Without considering the exogenous factors, time-series models seem to be the immediate fit to our task. Autoregressive integrated moving average (ARIMA) (Percival, Walden et al. 1993; Hamilton 2020) is arguably the most classical model for predicting time series data. The underlying ideas of ARIMA include auto-regression (i.e., the output variable is regressed on its own historical values), moving average (i.e., the regression error is a linear combination of error terms whose values occurred currently and at various times in the past), and integration (i.e., the data values have been replaced with the difference between their values and the



Figure 3: Distinct demand curves of different locations. The x-axis is the date, and the y-axis is the demand.

previous values). Each  $f_{\theta_i}(\cdot)$  can be represented as the following ARIMA(p,d,q) model (Nau 2020):

$$\hat{y}_{i,t} = \mu_i + \varphi_{i,1} \dot{y}_{i,t-1} + \dots + \varphi_{i,p} \dot{y}_{i,t-p} - \phi_{i,1} e_{i,t-1} - \dots - \phi_{i,q} e_{i,t-q}$$
(3)

Here  $\mu_i$  is a location specific constant, p is the number of autoregressive terms, and q is the number of lagged forecast errors in the prediction equation.  $\dot{y}_{i,t}$  is the  $d^{th}$  order difference of  $y_{i,t}$  to make it stationary. For example, when d = 1,  $\dot{y}_{i,t} = y_{i,t} - y_{i,t-1}$ .

Table 1: ARIMA model results

Location	1	2	3	4
Model	(0, 0, 0)	(0, 0, 0)	(2, 2, 4)	(0, 2, 4)
RMSE	7.112	4.284	5.714	3.177



Figure 4: Scatter plots of demand vs different factors. x-axis is the normalized value of the underlying factor, y-axis is demand.

Table 2: P-values of univariate linear regression. Features with p-values smaller than 0.05 are highlighted.

Feature	Solar radiation	Humidity	Temperature	Wind	Ferry	Library	Blue Bike
p-value	0	0	0	0.597	0	0.732	0.007
Feature	Snow cover	Snow depth	Snowfall	Snowmelt	Surface pressure	Precipitation	Cloud coverage
p-value	0.072	0.235	0.373	0.100	0.324	0.636	0.077

**Preliminary results** The prediction results of the ARIMA model in terms of RMSE are shown in Table 1. We can see that for locations 1 and 2, the best prediction is obtained by ARIMA(0,0,0), meaning that the fitted model is a constant value plus a white noise. This indicates that there is no significant trend in the demand time series for the two locations. For locations 3 and 4, the best results are respectively obtained by ARIMA(2,2,4) and ARIMA(0,2,4), indicating that there may be a certain trend that is captured by the model. Considering the scales of demand for different locations in Figure 3, the ARIMA-based predictions are reasonable for some locations (e.g., 3 and 4), but can be substantially improved for the other locations. We will revisit this in Section 5 where we discuss the final results.

# 4.3 Extract information from public data based on correlation analysis

Training a model solely from historical demand would miss contextual information about demand dynamics, where the exogenous factors are completely neglected. This could lead to arbitrarily bad predictions, especially when there is no significant trend of any order (e.g., for locations 1 and 2). After discussing with practitioners from The Family Van, we identify two important types of exogenous factors that can affect demand dynamics, i.e., weather and foot traffic.

It is intuitive that weather can affect the client demand. For example, people usually tend not to go out when the temperature is either too high or too low. To extract weather information, we use the public weather data from the (US) National Oceanic and Atmospheric Administration (NOAA)<sup>2</sup> and Copernicus.<sup>3</sup> These include weather attributes such as temperature, precipitation, solar radiation, humidity, snowfall, etc (see Table 2 for a complete list that we use).

In the meanwhile, we observe that the client demand is higher when there is more crowd around the vans. This indicates that *foot traffic* is an important factor of demand. Unfortunately, to the best of our knowledge, no foot traffic data in Boston are publicly available. Therefore, we use 3 other types of traffic information as surrogates of foot traffic, namely the public library usage data which is obtained from Analyze Boston,<sup>4</sup> the Massachusetts Bay Transportation Authority (MBTA) ferry usage data from Open Data MBTA,<sup>5</sup> and Blue Bike usage data from Bluebikes.<sup>6</sup>

Though we have a long list of features that could be potentially useful for prediction, it is often harmful to include any feature in a model, especially when a feature is not or very weakly correlated to the output. To measure the dependency of demand w.r.t. the features and therefore filter out informative features, we perform an Ordinary Least Squares (OLS) linear regression of demand vs each individual feature. Figure 4 shows the scatter plots of demand vs the feature value for some selected features. Each dot in the scatter plots represents a day's data point. The fitted linear curves are shown as the lines in the scatter plots. The p-values of the fitted univariate linear regression models are shown in Table 2. We then extract features whose associated p-values of the linear regression are smaller than 0.05. This returns 5 features: solar radiation, humidity, temperature, ferry usage and Blue Bike usage.<sup>7</sup> Surprisingly, 2 of the 3 features (Blue Bike usage ferry usage) that we consider as proxies for the foot traffic factor are highly correlated with demand.

## 4.4 Intermediate predictions

The previous section shows important insights on which features are informative of the mobile clinics demand. Though these features are desirable, the main barrier from actually using these features is that, the *future* values of these features are not available at *the time of prediction*. For example, when making predictions on Sunday, the blue bike usage information of the next Wednesday is not known a-priori.

For weather related features, fortunately, we can obtain reasonably accurate weather forecasts as the estimated fu-

<sup>&</sup>lt;sup>2</sup>https://www.noaa.gov/

<sup>&</sup>lt;sup>3</sup>https://www.copernicus.eu/en

<sup>&</sup>lt;sup>4</sup>https://data.boston.gov/

<sup>&</sup>lt;sup>5</sup>https://massdot.maps.arcgis.com/home/group.html?id= c5397b0d18d844c6a63b195a75ddf39b#overview

<sup>&</sup>lt;sup>6</sup>https://www.bluebikes.com/

<sup>&</sup>lt;sup>7</sup>Note that correlation does not imply causation. To analyze causal relations, a causal inference procedure is needed.



Figure 5: Curves of ground truth variable values (blue) and values based intermediate predictions (red).

ture values. However, there are no known tools or sources to obtain the estimated future values for the surrogate foot traffic features. To overcome this issue, instead of assuming the future values are available, we first use the time series models (e.g., ARIMA) to make *intermediate* predictions of the future values for these features, and then use the predicted values as input. Formally, the intermediate predictions are represented as:

$$\hat{x}_{i,t} = \mu_i^X + \varphi_{i,1}^X \dot{x}_{i,t-1} + \dots + \varphi_{i,p}^X \dot{x}_{i,t-p} - \phi_{i,1}^X e_{i,t-1}^X - \dots - \phi_{i,q}^X e_{i,t-q}^X,$$
(4)

where the superscript X is used to distinguish the notations for  $\hat{y}$  in Eq.(3).  $x_{i,t}$  is a generic representation of any features that are used as intermediate variables, and  $\dot{x}_{i,t}$  is the corresponding  $d^{th}$  order difference of  $x_{i,t}$  for stationarity.

The advantage of intermediate prediction vs an end-toend architecture where Eq.(4) is directly incorporated into Eq.(2) is two-fold. First, the model is trained with loss defined directly on the intermediate variable, instead of the final client demand. The latter introduces more noise. Second, we can use more data (rather than the target time period) to train the intermediate models.

Figure 5b shows the predicted and ground truth curves for the 5 selected features. We can see that the predictions are overall reasonably capturing the trends of the ground truth curves. For temperature, solar radiation and humidity, the intermediate predictions (obtained from weather forecast) are highly aligned with the ground truth values. This demonstrates that the weather forecasts are reliable. The same holds for ferry usage, where the prediction curve is also highly aligned with the ground truth. For bike usage intermediate prediction, there is a noticeable bias due to the discrepancy between the training and test data. Nonetheless, the trend is still largely captured by the prediction.

## 4.5 Integrating with different ML models

With the intermediate predictions, we show how we integrate these intermediate variables with different types of ML models to make the demand prediction. For non-time-series (NTS) models, this means

$$\hat{y}_{i,t} = f_{\theta}(\hat{x}_{i,t}; t), \tag{5}$$

where  $\hat{x}_{i,t}$  is the estimated feature values based on the intermediate predictions in Eq.(4). Note that we still feed the time information t to this model to compensate for the loss of temporal information. Alternatively, we have also used recurrent neural networks which combine the historical demand and (estimated) contextual features:

$$\hat{y}_{i,t} = f_{\theta}(y_{i,1}, \dots, y_{i,t-1}; \hat{x}_{i,t}).$$
(6)

### 4.6 Separating predictions for repeat and non-repeat clients

Another important observation, as shown in Fig. 6, is that the demand patterns of the repeat and non-repeat clients for different locations are distinct. Take location 2 as an example, the demand curve of the non-repeat clients has a significantly larger fluctuation over time, whereas the demand of the repeat clients is more stationary. Inspired by this, we propose an alternative training method, where we first train 2 separate models for each group of the repeat and non-repeat clients, and then sum the predictions of the two models as the final prediction. This essentially means that we will train a model for each of  $\hat{y}_{i,t}^R$  (repeat client demand) and  $\hat{y}_{i,t}^N$  (nonrepeat client demand) in either Eq.(5) or Eq.(6).



Figure 6: Demand curves of repeat vs non-repeat clients

## 5 Main results

**Experiment settings** Our demand dataset from The Family Van spans from July 2019 to March 2020. It contains daily demand data for 4 locations in Boston. The Family Van operates on 4 days a week, one day per location. Therefore, we have one aggregated data point per week per location during the target time period, totalling 20-30 data points per location and 110 data points for all locations. We use the first 80% of the data to train a model for each location and evaluate it on the rest 20% of data. In practice,

row number	method $\setminus$ location	Location 1	Location 2	Location 3	Location 4
1	Historical average	7.06	4.20	8.91	4.45
2	ARIMA	7.11	4.28	5.71	3.18
3	NTS-AB	$6.14\pm0.0000$	$3.95\pm0.0000$	$6.47\pm0.0000$	$3.17\pm0.3776$
4	NTS-Ab	$6.17\pm0.0000$	$\textbf{3.09} \pm 0.5938$	$6.48\pm0.0000$	$\textbf{2.90} \pm 0.0000$
5	NTS-aB	$3.40 \pm 0.0000$	$3.77\pm0.4502$	$6.80\pm0.0000$	$3.39\pm0.0000$
6	NTS-ab	$3.65\pm0.0000$	$3.93\pm0.0000$	$6.69\pm0.0000$	$3.31 \pm 0.0199$
7	RNN-AB	$7.96 \pm 1.8481$	$10.31 \pm 1.1406$	$6.03\pm0.2806$	$4.96\pm0.3454$
8	RNN-Ab	$7.64 \pm 0.7485$	$7.39 \pm 1.4918$	$6.66\pm0.2225$	$6.75 \pm 1.0613$
9	RNN-aB	$8.33 \pm 0.7146$	$10.28 \pm 1.2289$	$7.38\pm0.6014$	$5.31 \pm 0.3027$
10	RNN-ab	$7.47\pm0.5451$	$7.11 \pm 1.8134$	$6.53\pm0.1624$	$5.63 \pm 0.4278$

Table 3: Demand prediction using our proposed techniques. For the set of NTS and RNN models, the result of the best model is given. The best result for each location is highlighted.

the set of NTS models we implement are: Linear regression, Ridge regression, Lasso regression, Lasso LARS, Tweedie regression, SGD, Logistic regression, MLP, Adaboost, Decision Tree (these are implemented via scikit-learn<sup>8</sup>) as well as XGBoost (Chen et al. 2015). For RNNs, we use 3 structures: the vanilla RNN, LSTM (Hochreiter and Schmidhuber 1997) and GRU (Cho et al. 2014). For these models, we run 20 times of training, and report the average RMSE and standard deviation of RMSE of the best NTS or RNN model, respectively.

For intermediate features, we can either train with ground truth feature values, and test with predicted feature values (with Eq.(4)), or both train/test with predicted feature values. We use **A** to denote the former setting and **a** to denote the latter. Similarly, we can either train a single model for both types of the repeat and non-repeat clients, or train 2 separate models. We denote them as **B** and **b**, respectively. Therefore, a training method **Ab** means we train with ground truth feature values, and test with predicted feature values, and at the same time train 2 separate models for the repeat/non-repeat clients. The main results are shown in Table 3. We will explore the answers to the questions in the following.

Are predictions accurate in general? Comparing the best result and the historical average in row 1 (which is the current practice of The Family Van), we can see that the RMSE respectively decreases by 51.8%, 26.4%, 35.9%, 34.8% for the 4 locations. Combing the overall scales of the demand for the 4 locations (see Figure 3), this demonstrates that *our proposed method is a reasonably reliable tool for demand prediction*.

What is the best ML model (if there is any)? Most surprisingly, RNN models are dominated by the NTS and ARIMA models. Our understanding is that deep models have more parameters, and thus tend to overfit in small datasets like ours. Second, NTS models are significantly better than ARIMA, except for location 3. Our understanding is that the demand curve (see Figure 3) for location 3 is noticeably declining, whereas this is not obvious for other locations, especially 1 and 2. This is also aligned with the observations of the best ARIMA configuration in Table 1. *There is no uniformly the best model for all the locations. The time series model is the best fit for location 3, whereas the NTS models are the best fit for the other locations.* 

Are exogenous factors and the intermediate predictions helping? Comparing the best results for each location in Table 3 and the results obtained by the pure time-series ARIMA model in Table 1 (the same as row 2 in Table 3), we can see that there are substantial improvements for 3 out of the 4 locations. Notably, for location 1, the RMSE of the ARIMA model is more than twice of the feature-based model. ARIMA performs slightly better in location 3. This shows that exogenous factors are critical to the demand prediction. Another interesting observation is that, when both training and testing on the intermediate feature values, there is a substantial gain in accuracy for location 1 (comparing row 3 vs row 5 or row 4 vs row 6). Our hypothesis is that the intermediate prediction of features like Blue Bike usage is biased. Therefore, when training on observed features and testing on predicted features, there is a gap between the two values. This gap is somehow removed when both training and testing on the intermediate feature values.

Is separating predictions for repeat and non-repeat clients helping? Comparing row 3 vs row 4, we can see that separating predictions leads to substantial improvements in locations 2 (3.95 to 3.09) and location 4 (3.17 to 2.9), while maintaining very close results in the other 2 locations. Our hypothesis for the improvement is that for locations 2 and 4, the demand patterns of the repeat and non-repeat clients are more distinct, as shown in Figure 6. Comparing row 5 vs row 6, we can see that the results are close for both training methods in all of the locations (separate prediction is slightly better in locations 3 and 4, and slightly worse in locations 1 and 2). This shows that *separate prediction helps in improving the overall prediction accuracy*.

## 6 Conclusion

As far as we know, we are the first to explore weekly/daily dynamic demand prediction of mobile clinics using AI. We propose a novel learning framework that uses publicly available data for prediction, together with multiple innovations

<sup>&</sup>lt;sup>8</sup>https://scikit-learn.org/stable/

that are customized into solving the demand prediction problem. Empirical results on real-world datasets from The Family Van demonstrate that our proposed approach has substantial improvement in accuracy compared to the experiencebased estimation. Our study provides a brand-new angle to mobile clinics demand prediction with completely public data, which has a huge potential impact when broadly deployed. As a future work, we are actively exploring how our prediction algorithm can be deployed to help The Family Van's daily scheduling of staff and healthcare resources.

## 7 Acknowledgments

We thank Ariel Procaccia, Han-Ching Ou, and Herman Saksono for the early discussions and brain storming. We thank Karen Li for early exploratory analysis and Rainelle Walker-White for feedback on the problem background.

Chen and Biswas were supported by the Harvard Center for Research on Computation and Society.

### References

Bhowmik, T.; and Eluru, N. 2021. A Comprehensive County Level Framework to Identify Factors Affecting Hospital Capacity and Predict Future Hospital Demand. *medRxiv*.

Bolen, S. D.; Sage, P.; Perzynski, A. T.; and Stange, K. C. 2016. No moment wasted: the primary-care visit for adults with diabetes and low socio-economic status. *Prim. Health Care Res. Dev.*, 17(1): 18–32.

Chen, T.; He, T.; Benesty, M.; Khotilovich, V.; Tang, Y.; Cho, H.; et al. 2015. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4): 1–4.

Cho, K.; van Merriënboer, B.; Gulcehre, C.; Bahdanau, D.; Bougares, F.; Schwenk, H.; and Bengio, Y. 2014. Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. In *Proceedings of the* 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), 1724–1734.

Côté, M. J.; and Smith, M. A. 2018. Forecasting the demand for radiology services. *Health Systems*, 7(2): 79–88.

de Peralta, A. M.; Gillispie, M.; Mobley, C.; and Gibson, L. M. 2019. It's All About Trust and Respect: Cultural Competence and Cultural Humility in Mobile Health Clinic Services for Underserved Minority Populations. *J. Health Care Poor Underserved*, 30(3): 1103–1118.

Drackley, A.; Newbold, K. B.; Paez, A.; and Heddle, N. 2012. Forecasting Ontario's blood supply and demand. *Transfusion*, 52(2): 366–374.

Hamilton, J. D. 2020. *Time series analysis*. Princeton university press.

Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation*, 9(8): 1735–1780.

Khaldi, R.; El Afia, A.; Chiheb, R.; and Faizi, R. 2017. Artificial neural network based approach for blood demand forecasting: Fez transfusion blood center case study. In *Proceedings of the 2nd international Conference on Big Data, Cloud and Applications*, 1–6. Mackey, K.; Ayers, C. K.; Kondo, K. K.; Saha, S.; Advani, S. M.; Young, S.; Spencer, H.; Rusek, M.; Anderson, J.; Veazie, S.; Smith, M.; and Kansagara, D. 2021. Racial and Ethnic Disparities in COVID-19-Related Infections, Hospitalizations, and Deaths : A Systematic Review. *Ann. Intern. Med.*, 174(3): 362–373.

Majeed, B.; Peng, J.; Li, A.; Lin, Y.; and Delgado, R. I. 2021. Forecasting the demand of mobile clinic services at vulnerable communities based on integrated multi-source data. *IISE Transactions on Healthcare Systems Engineering*, 11(2): 113–127.

Nau, R. 2020. Introduction to ARIMA models. https://people.duke.edu/~rnau/411arim.htm#arima010.

Odlum, M.; Moise, N.; Kronish, I. M.; Broadwell, P.; Alcántara, C.; Davis, N. J.; Cheung, Y. K. K.; Perotte, A.; and Yoon, S. 2020. Trends in Poor Health Indicators Among Black and Hispanic Middle-aged and Older Adults in the United States, 1999-2018. *JAMA Netw Open*, 3(11): e2025134.

Percival, D. B.; Walden, A. T.; et al. 1993. *Spectral analysis for physical applications*. cambridge university press.

Qian, D.; Pong, R. W.; Yin, A.; Nagarajan, K.; and Meng, Q. 2009. Determinants of health care demand in poor, rural China: the case of Gansu Province. *Health Policy and Planning*, 24(5): 324–334.

Reed, C.; Rabito, F. A.; Werthmann, D.; Smith, S.; and Carlson, J. C. 2019. Factors associated with using alternative sources of primary care: a cross-sectional study. *BMC health services research*, 19(1): 1–9.

Stephanie, W.; Hill, C.; Ricks, M. L.; Bennet, J.; and Oriol, N. E. 2017. The scope and impact of mobile health clinics in the United States: a literature review. *International journal for equity in health*, 16(1): 1–12.

Syed, S. T.; Gerber, B. S.; and Sharp, L. K. 2013. Traveling towards disease: transportation barriers to health care access. *J. Community Health*, 38(5): 976–993.

Vance, K.; Howe, W.; and Dellavalle, R. P. 2009. Social internet sites as a source of public health information. *Dermatol. Clin.*, 27(2): 133–6, vi.

Woolliscroft, J. O. 2020. Innovation in Response to the COVID-19 Pandemic Crisis. *Acad. Med.*, 95(8): 1140–1142.

Zimmerman, F. J.; and Anderson, N. W. 2019. Trends in Health Equity in the United States by Race/Ethnicity, Sex, and Income, 1993-2017. *JAMA Netw Open*, 2(6): e196386.