

1 **Using Neural Network for Predicting Hourly Origin-Destination Matrices**
2 **from Trip Data and Environmental Information**

3
4 Ehsan Hassanzadeh¹, Zahra Amini^{*,1}

5
6 ¹ Department of Civil Engineering, Sharif University of Technology, Tehran, Iran

7
8 E-mail Addresses: ehsan.hassanzadeh@student.sharif.edu (E. Hassanzadeh),
9 zahra.amini@sharif.edu (Z. Amini)

10
11 Submitted for review and possible publication in
12 Scientia Iranica

13
14
15 Corresponding author:

16 Zahra Amini
17 Azadi Avenue
18 Tehran
19 Iran

20 Email: zahra.amini@sharif.edu

21 Phone:
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 **Biography**

46 **Ehsan Hassanzadeh** is a graduated Master's student in transportation engineering at the
47 Sharif University of Technology, Tehran, Iran. In his master's studies, he has focused on
48 transportation network modeling and traffic percolation under the supervision of Dr. Zahra
49 Amini. During his Master's career, he has also focused on applications of Machine Learning
50 methods in Transportation Engineering, specifically on the application of Neural Networks in
51 estimating travel demands. He holds a bachelor's degree in Civil Engineering from the Ferdowsi
52 University of Mashhad and he has started his Ph.D. studies in Transportation Engineering at the
53 University of British Columbia (UBC).

54 **Zahra Amini** is currently an Assistant Professor at the Department of Civil Engineering,
55 Sharif University of Technology. She completed her bachelor's degree in 2014 and her master's
56 degree in 2015, in Civil Engineering at University of California, Berkeley. She obtained her
57 Ph.D. in Highway and Traffic Engineering, in 2018 at University of California, Berkeley. Her
58 research interests are Intelligent Transportation System (ITS), traffic theory and control
59 strategies, and transportation system operation and management.
60

61 **Abstract**

62 Predicting Origin-Destination demand has always been a challenging problem in transportation.
63 Conventional demand prediction methods mainly propose procedures for forecasting aggregated temporal
64 Origin-Destination (OD) flows. In other words, they are primarily unable to predict short-term demands.
65 Another limitation of these models is that they do not consider the impact of environmental conditions on
66 trip patterns. Furthermore, OD demand prediction requires two individual steps of modeling: trip
67 generation and trip distribution. This article presents a framework for predicting hourly OD flows using
68 the Neural Network. The proposed method utilizes trip patterns and environmental conditions for
69 predicting demands in single-step modeling. A case study on New York City Green Taxi 2018 trip data is
70 done to evaluate the method, and the results demonstrate that the network has reasonably accurate OD
71 flows predictions.
72

73 **Keywords:** O/D demand prediction; Short-term prediction; Neural Network; Machine
74 Learning, Trip Generation
75

76 **Highlights**

- 77
- 78 • Short-time O/D flow prediction is proposed to be obtained by Neural Network models
 - 79 • The proposed Neural Network model uses past trip and environmental condition data
 - 80 • The proposed model may replace aggregate distribution models for short-time predictions

81 **1. Introduction**

82 In recent years, demand for public transport has increased significantly due to urban
83 development and population growth. One option to meet this demand is public transportation
84 network expansion, which is expensive and has many limitations [1]. A more appropriate
85 solution is network management with the available facilities. Network management includes
86 strategies and policies to fulfill demand in the system and utilize the facilities more effectively
87 [2]. A preliminary requirement if devising such management plans is predicting and modeling
88 users' travel behavior.

89 One of the most popular models in demand modeling is the Four-Step Model (FSM), which
90 contains the following steps: 1. trip generation, 2. trip distribution, 3. mode choice, and 4. route
91 choice [3]. Although this paper mainly focuses on the first two steps and assumes predicting O/D
92 demands for a single mode, some dynamic factors are considered that could alter users' mode
93 choices. The second step, trip distribution, distributes generated trips to match destinations. In
94 the third step, mode choice and trip modes proportions are specified, and the user behavior is
95 modeled using the Consumer Choice theory. Based on this theory, consumers' preference is
96 affected by utility functions, which are not deterministic [4–6]. By assuming a uniform
97 distribution for the random component of the utility function across the modes, it can be
98 concluded that the difference in utilities is only due to the difference in the systematic part.
99 However, it should be noted that researchers have assumed that the random components are
100 independent but non-identically distributed [4].

101 Nonetheless, reviewing the random component's distributions is beyond the scope of this
102 paper, and since it will not be used in this research's modeling process, the distribution is
103 assumed to be uniform to explain the impact of the systematic component. As a result, users'
104 mode choice behavior can be modeled and predicted by considering systematic utility
105 components, including consumers' socio-economic attributes, the vehicle's operational
106 characteristics, and the trip distribution table [4]. Moreover, previous studies show that
107 environmental conditions (e.g., weather data, land use, and other related parameters) have
108 substantial impacts on travel times, public transportation systems quality, and daily travel
109 behaviors [7–10]. As a result, these parameters can be considered in the consumer's utility
110 function. The relationship of these parameters with the utility function (e.g., linear or nonlinear)
111 would be determined by the Neural Network in this project by creating dummy variables in the
112 hidden layers; hence, this paper will not discuss the possible relationships with the utility
113 function. Given that the other parameters, such as the socio-economic attributes (e.g., the impact
114 of income on the mode choice), do not vary over short periods, this paper disregards such
115 parameters as they are assumed to be unchanged over the study period.

116 This paper aims to predict hourly OD flows for a single specific mode using the Neural
117 Network (NN) without users' information (e.g., income, car ownership). To predict OD flows,
118 input parameters reflecting the consumer's utility and other parameters regarding trip information
119 are used. Then to evaluate the proposed method, this study uses New York City Green Taxi 2018
120 trip data and New York City 2018 weather data. The trip data consists of 8.81 million trip
121 information, including trips' origin and destination zones number, trip distance, and other related
122 trip information. Another dataset that is used for training the model includes hourly weather data
123 of the desired location.

124 The rest of this article is organized as follows. Section 2 reviews related works on OD
125 prediction and studies using a similar dataset. Section 3 discusses the proposed framework to
126 predict OD flows in detail. Then, section 4 describes datasets used to evaluate the proposed
127 model and data verification by investigating existing trip patterns. Section 5 discusses model
128 results, and the final section provides the conclusion.

129 130 **2. Literature review**

131 Demand modeling has been a prominent research area in transportation for years, and the
132 FSM has been one of the most comprehensive approaches for demand modeling. This approach
133 underlies methods to predict mode-specific demand [11]. The Gravity model is widely used for
134 trip distribution in the trip generation step of the FSM. This model distributes trips between

135 zones based on the zones' relative attraction and a function of distances between zones [12]. The
136 model is calibrated on a single OD table, including aggregated trip data. Thus, applying the
137 Gravity model for estimating hourly OD flows is not practical. Moreover, the Gravity model
138 considers a limited number of parameters to calibrate the model. Studies have shown that the
139 Gravity model output has little similarity with the observed data [13]. It can be inferred from the
140 Gravity model that the origin and destination zones and a function of distances between the
141 zones (either temporal or spatial) should be considered in the modeling.

142 The third step of the FSM involves using discrete choice models to understand user behavior
143 when selecting transportation modes. The Logit and Probit models are commonly used in this
144 step, which use utility functions to determine choice probabilities [14,15]. However, collecting
145 user information to define these utility functions can be challenging. These models have
146 limitations, such as assuming a constant relationship between dependent and independent
147 variables, making them inflexible and unable to adapt dynamically. They also perform poorly
148 when input variables are multicollinear [16–18]. This paper focuses on parameters that can affect
149 user choices based on environmental conditions, rather than explicitly modeling mode choice.

150 Recent studies have explored the use of Neural Networks (NN) as an alternative to
151 traditional methods for predicting OD matrices and mode choices [19–23]. Researchers have
152 compared the accuracy and performance of NN models with other statistical methods, such as
153 the Multinomial Logit Model (MNL) [24], mode choice modeling [25], and Bayesian Model
154 Tree [26]. Xiong et al. [27], proposed a framework that used Graph Neural Networks (GNN) and
155 Kalman filters to predict OD flows based on historical link flows. Yaldi et al. [28] used NN
156 models with three input parameters to predict trip flows. However, these approaches have
157 limitations in considering the factors impacting OD flows and user behavior. The current paper
158 proposes a new approach that uses NN models to predict trip flows based on trip patterns,
159 environmental conditions, and consumer preferences. Like the method proposed by Xiong et al.
160 [27], the framework used by Yaldi et al. [28] limits input parameters to trip interchange
161 attributes, ignoring environmental attributes affecting users' behaviors. In contrast, the current
162 paper implements the NN to predict trip flows using trip patterns and environmental conditions
163 considering consumer preferences.

164 Researchers have realized that environmental conditions may impact traffic patterns in
165 various ways. Liu et al. claimed that weather parameters, including temperature, snowfall, and
166 precipitation, substantially impact travel behaviors [7]. They showed that these weather
167 parameters affect all travel modes, including pedestrian walking, bicycle, private car, and public
168 transport. In another study, Rudloff et al. evaluated relations between weather conditions and trip
169 patterns using mode choice models. They estimated choice models' parameters based on
170 household survey data from Vienna, employing the maximum likelihood approach. Their results
171 showed that weather conditions significantly influence transport choice and travel behavior [8].
172 Hyland et al. investigated the effects of weather conditions on travel mode choice using a stated
173 preference (SP) survey in Chicago and realized that commute choice patterns differ vastly in
174 various weather conditions. Furthermore, they claimed that the impacts of weather on mode
175 choices vary across the community [29]. Thus, the present paper considers weather conditions as
176 effective environmental parameters while training the network for predicting OD flows.

177 The NN is mainly trained on the existing trip patterns to learn future predictions.
178 Consequently, it is essential to create parameters considering different trip patterns to have a
179 more accurate estimate of the future. Studies have shown that different trip behaviors are
180 observed on weekends and holidays compared to workdays. Dong et al. [30] used trajectory data

181 collected from ride-hailing services in Beijing, China, to investigate urban trip patterns. Their
182 results showed tangible differences in trip distributions between particular zones. Specific hourly
183 patterns also justify considering the hour of the day as an effective parameter in network training.
184 They observed a notable difference between workday and non-workday trip patterns for various
185 trip purposes. Other researchers have also shown the importance of the time of the day in
186 predicting OD flows [31–33]. These results reveal the importance of considering weekends,
187 holidays, and hour of the day in model training.

188 Another aspect of this research is dealing with big data for transportation analysis. In this
189 regard, numerous research uses big data for various types of analysis. To name a few related
190 research, [34,35] use big data for analyzing a specific mode of transportation. The latter also
191 focuses on the impacts of COVID-19 on bike-sharing systems. Another similar approach to
192 dynamically predict trip patterns using the NN is to apply agent-based day-to-day models. In this
193 field, many papers focus on trip-related information and how it can impact traffic conditions by
194 applying agent-based models [36,37]. Based on the nature of the problem, which includes
195 various parameters impacting traffic patterns and mode choices, this paper opts to utilize the NN
196 for predicting OD flows.

197 This article uses 2018 New York City taxi data from NYC Open Data to evaluate the
198 proposed framework's performance. Related works on similar datasets are as follows. Deri et al.
199 used similar 2010-2013 New York City taxi data and presented a solution for estimating taxi
200 trajectories using Dijkstra's algorithm with a significantly reduced computation time [38]. In
201 another study, Freire et al. discussed cleaning Spatio-temporal data. They used 2008-2012 New
202 York City taxi data to observe the anomalies in the dataset. Results showed that data exploration
203 needs users' assistance, and the lack of adequate information about events prevents the system
204 from discerning anomalies [39]. Patel et al. proposed an approach to visually explore big OD
205 data and determine average hourly drivers' revenue. They used 2014 New York City taxi data to
206 evaluate their method. Unlike related works using a similar dataset [40], this paper aims to
207 predict OD flows considering the abovementioned parameters.

208 209 **3. Materials and Methods**

210 This section describes data cleaning procedures and obtaining various input and output
211 parameters required for modeling and the network structure. Input parameters (independent
212 parameters in modeling) indicate parameters used as the network's input to predict output
213 parameters (dependent parameter in modeling), the OD flow per hour. In other words, hourly
214 input data are used to predict the hourly OD flows. This section divides the proposed algorithm's
215 procedure into four major steps, as summarized below:

- 216
217 • Step 1: Obtaining Input Parameters (Independent Parameters) data for training and testing the NN
218 (See Section 3.1). This step consists of the following minor steps:
 - 219 ○ 1A: Obtaining OD zone IDs and hour of the day (See Section 3.1.1)
 - 220 ○ 1B: Obtaining interzonal travel times (See Section 3.1.2)
 - 221 ○ 1C: Obtaining binary parameters (See Section 3.1.3)
- 222 • Step 2: Obtaining output parameters data (number of trips for each pair of OD zones at each time
223 step) for training and testing the NN (See Section 3.2).
- 224 • Step 3: Cleaning obtained data to remove any outliers that may deteriorate the NN accuracy (See
225 Section 3.3).
- 226 • Step 4: Reshaping input and output data matrices to be fed into the NN for training (See Section
227 3.4).

- *Step 5*: Standardize data to avoid bias in training the NN (See Section 3.5).
- *Step 6*: Building the NN's structure and training the network (See Section 3.6).

According to the possible factors affecting trips described in earlier sections, this paper considers the network's input parameters as follows:

- Interzonal travel times, including calculated hourly travel times for all possible OD pairs
- Origin zone, defined by a unique ID
- Destination zone, defined by a unique ID
- Hour of the day, specified by a number within the range zero to 23
- Weekend/ weekday binary classification
- Holiday/ nonholiday binary classification
- Temperature, including hourly resolution records
- Precipitation, including hourly resolution records
- Snow depth, including hourly resolution records

This paper assumes that certain parameters such as passenger count and fare amount cannot be determined without access to corresponding demand data or algorithms used to calculate these parameters. Therefore, these parameters are not used as inputs in the model. Additionally, input parameters with intercorrelation are omitted, such as distance between origin and destination zones. The procedures for obtaining each parameter are explained in detail below.

3.1. Obtaining input data

As described earlier, each hourly input record comprises nine factors. Three of these nine parameters, which are temperature, precipitation, and snow depth, can be obtained directly from datasets for each time step. Obtaining the remaining six requires additional steps, which are described below.

3.1.1. Obtaining OD zone IDs and hour of the day

In this research, the "hour of the day" variable is defined as the departure hour for each trip. This definition may raise an error since trips are not necessarily finished in the same hour as they started. Since it is assumed that the desired trips only consist of urban trips, trip durations would be reasonably short; hence the error is negligible. Zone IDs can be defined as the assigned IDs for each Traffic Analysis Zones (TAZs). Therefore, each trip's destination and origin can be determined with two IDs demonstrating its origin and destination. Traffic Analysis Zones can be specified using the available datasets for the research area or by defining the TAZs using the available methods [41].

3.1.2. Obtaining interzonal travel times

The Gravity model examines how distances between zones impact OD flows, but this article suggests using hourly travel times between zones to consider the impact of traffic flows. Hourly travel times can be obtained from various services like Google Maps, but the calculation of shortest paths requires real-time traffic data. Since this research aims to predict hourly OD flows, all input and output data should be aggregated into hourly records. The article proposes calculating average travel times between OD pairs in each hour after removing outliers to represent the hourly travel time for all the trips between the OD zones. Section 3.3 provides more

272 details on the removal of outliers. Then, interzonal hourly travel times are obtained as a linear
 273 matrix, TT_k , according to Equation (1).

$$274$$

$$275 \quad TT_k = [tt_{1,1,k}, \dots, tt_{1,n,k}, \dots, tt_{i,j,k}, \dots, tt_{n,n,k}] \quad (1)$$

276

277 TT_k = linear travel time matrix in the hour k , $tt_{i,j,k}$ = travel time for the ij OD pair for the hour k .

278

279 Depending on the dataset used, travel times for the OD pairs may be obtained using a
 280 particular method available. As will be discussed in Section 4.1, the dataset used in this study
 281 includes the start time and end time for each trip record. Therefore, each trip's travel time can be
 282 simply calculated by computing the in-time vehicle for each trip.

283

284 3.1.3. Obtaining binary parameters

285 As described in section 2, users have different traffic behaviors on weekends and holidays
 286 than on regular weekdays. Two binary parameters are defined to address this variation:
 287 "Weekend" and "Holidays," parameters which indicate whether the trip was on the weekend or
 288 holiday or not, respectively. These parameters' values are equal to zero if the desired day is not a
 289 holiday or a weekend. It is worth mentioning that the holidays can be specified using the national
 290 holidays' list for the desired database. After obtaining all parameters as discussed, linear matrices
 291 of hourly attributes, att_k can be created according to Equation (2). The data source for obtaining
 292 these parameters will be discussed in Section 4.2.

$$293$$

$$294 \quad att_k = [hr_k, weekend_k, holiday_k, temp_k, PCP_k, SD_k]$$

295 (2)

296 att_k = linear attributes matrix of the hour k , and for the hour k : hr_k = the hour k of the day, $Temp_k$ =
 297 the hour k hourly temperature, PCP_k = the hour k hourly precipitation, SD_k = the hour k snow depth.

298

299 3.2. Obtaining output data

300 The network's output parameter for each hour is an OD flow matrix showing trip counts
 301 between each OD pair. The OD flow matrices at each hour are created by counting the trips
 302 between each pair of OD in the trips database. In this study, origins and destinations are
 303 considered Traffic Analysis Zones (TAZs), which can be specified by assigning zone IDs. The
 304 hourly OD matrix is created as shown in Equation (3).

$$305$$

$$306$$

$$307 \quad OD_k = \begin{bmatrix} T_{1,1} & \cdots & T_{1,n} \\ \vdots & T_{i,j} & \vdots \\ T_{n,1} & \cdots & T_{n,n} \end{bmatrix} \quad (3)$$

308

309 $T_{i,j,k}$ = trip counts from zone i to zone j , for the hour k . k = the data record index representing the hour
 310 k

311

312
 313 It should be noted that hourly trips are counted based on their start time (i.e., departure time).
 314 These matrices are then reshaped to linear matrices, as shown in Equation (4), to simplify the
 315 network's training process since it would be less baffling to acquire one row of data per record
 316 when feeding the input data to the network.

$$317$$

$$318$$

$$319 \quad OD_k = [T_{1,1,k} \dots T_{1,n,k} \quad T_{i,j,k} \dots T_{n,n,k}] \quad (4)$$

320
 321
 322 These matrices are then added to the final output matrix, T , according to the occurrence time,
 323 starting from the first hour of the initial day ($k = 0$) to the last hour of the last day in the period (k
 324 $= t$). Thus, the output matrix, T , would be a $(t, n \times n)$ dimensional matrix according to Equation
 325 (5). Each row of the output matrix (i.e., dependent variable) indicates hourly trip counts for an
 326 OD pair. The predicted results after training the model will also follow the same format.

$$327$$

$$328 \quad T = \begin{bmatrix} OD_0 \\ \vdots \\ OD_t \end{bmatrix} \quad (5)$$

329 330 331 **3.3. Data Cleaning**

332 The existence of errors in data will result in bias in the network's training process. As a
 333 result, possible errors should be omitted from the data before using it for training. This section
 334 discusses removing outliers and possible errors from input and output data. To do so, the Z-score
 335 is calculated for each record to identify outliers in the data. The Z-score indicates the distance
 336 between the observed value and the sample's mean in the standard deviation units [42]. The Z-
 337 score can be calculated using Equation (6).

$$338$$

$$339 \quad z = \frac{x - \mu}{\sigma} \quad (6)$$

340
 341
 342 z = the standard score, x = the observed value, μ = the mean of the sample, σ = the standard
 343 deviation of the sample

344
 345 After calculating the Z-score, records with $|z| \geq 3$ are considered outliers (preserving 99.8% of
 346 the data range). It should be noted that the threshold for removing the records is calculated after
 347 investigating the results by examining different thresholds. This procedure should be done for all
 348 input and output parameters with possible errors. Besides, constraints should be set for each
 349 parameter to ensure that all remaining data are valid. For example, travel time values should be
 350 positive and over 60 seconds. Values exceeding these ranges should be omitted based on the
 351 parameter range.

352
353

3.4. Reshaping data matrices

354 So far, the origin and destination zone IDs have not been determined in the input and output
355 matrices. As discussed before, each pair of OD has a unique traffic pattern. So, it is crucial to
356 consider origin and destination zones as parameters for training the model. Then, input and
357 output data are reshaped so that each row of data matrices represents dependent and independent
358 values for a specific OD pair in a specific hour. The final output matrix, T , would be as shown in
359 Equation (7).

$$360 \quad T = \begin{bmatrix} OD_0 \\ \vdots \\ OD_k \\ \vdots \\ OD_t \end{bmatrix} \rightarrow \begin{bmatrix} T_{1,1,0} \\ \vdots \\ T_{n,n,0} \\ \vdots \\ T_{i,j,k} \\ \vdots \\ T_{n,n,t} \end{bmatrix} \quad (7)$$

361
362
363
364
365
366

To create the final input matrix, X , the travel times matrices, TT_k , should be first reshaped similar to the output matrix. Then, trip attributes parameters can be appended, including origin and destination zone IDs for each record, duplicating common attributes for all the trips in the desired hour. As a result, the final input matrix, X , is according to Equation (8).

$$367 \quad X = \begin{bmatrix} tt_{1,1,0} & i & j & hour_0 & weekend_0 & holiday_0 & temp_0 & PCP_0 & SD_0 \\ & & & & & \vdots & & & \\ tt_{i,j,k} & i & j & hour_k & weekend_k & holiday_k & temp_k & PCP_k & SD_k \\ & & & & & \vdots & & & \\ tt_{i,j,t} & i & j & hour_t & weekend_t & holiday_t & temp_t & PCP_t & SD_t \end{bmatrix} \quad (8)$$

368
369
370
371
372

$tt_{i,j,k}$ = calculated travel time for the origin zone i and the destination zone j in the hour k
 i, j = origin and destination zone IDs for each record

3.5. Data Standardization

373 Due to the significant variances between parameter values (either the difference between
374 values of one parameter or the diversity between the data range of various parameters), the
375 network's training process may be biased. As a result, large trip counts in the OD matrix, which
376 are vital for modeling, could be recognized as outliers. All parameters' values are standardized in
377 their category to address this issue, making the mean of each parameter zero and the standard
378 deviation of parameter one. Standard values are calculated using Equation (9) [43].

379

$$380 x' = \frac{x - \mu}{\sigma} \tag{9}$$

381

382

383 x' = the standardized value, x = the observed value, μ = the mean of the sample, σ = the standard
384 deviation of the sample

385

386 It should be noted that the network's predicted data will be calculated in the normalized format
387 and must be converted to the original format for evaluating the model.

388

389 **3.6. The Neural Network's structure**

390 The Neural Network is a supervised machine learning method in which the network is
391 trained first using a set of data with pre-defined outputs. The network tries to minimize the
392 defined objective function to achieve the most desirable results by finding connections between
393 input and output nodes. The NN is composed of multiple layers including input, hidden, and
394 output layers. Hidden layers identify possible relations between parameters and provide a
395 representation of data with multiple layers of abstraction. Each layer has a specific activation
396 function to transmit the data format for the next layer. The NN optimizer updates weights based
397 on the gradients computed in each iteration through an iterative backpropagation process [44].

398 The NN used in this paper has two hidden layers, as illustrated in Figure 1. The input layer is
399 provided with the input parameters and transmits input data directly to the next layer via the
400 neurons. As mentioned earlier, the number of nodes in this layer equals the number of input
401 parameters, which is nine. The input and output layers' dimensions are 9 and 1, implying the
402 input feature vector's dimensionality and prediction value. The predicted value here is the hourly
403 OD flow, and each predicted value defines the predicted flow for a specific OD pair in a specific
404 hour of the desired period. To determine the number of nodes in the hidden layer, different
405 numbers of nodes can be chosen, and then the output results of the network can be compared. As
406 a general experimental rule, the number of nodes in the hidden layers is chosen close to the
407 average input and output number of nodes. After investigating the results with different counts of
408 nodes for the hidden layers, 7 and 5 nodes are finally considered for the hidden layers,
409 respectively, as shown in Figure 1.

410 Activation functions are added to the NN to convert the previous layer's output values into
411 desirable input values for the next layer. The ReLU (Rectified Linear Unit) activation function is
412 used for hidden layers in this network as it offers better performance and generalization than the
413 other activation functions used for predicting a numerical value [45]. The ReLU function can be
414 written as Equation (10).

415

$$416 ReLU(x) = \max(0, x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{if } x < 0 \end{cases} \tag{10}$$

417

418

419 As illustrated in Figure 1, the output layer connects the last hidden layer to the output values.
420 This layer has only one node, which is the network output value. Activation functions like the
421 Sigmoid function result in output values between 0 and 1 that predict categorical values. The

422 output values in this work are numerical; thus, the Linear function is used as an activation
423 function for the output layer, expressed as Equation (11).

$$424$$
$$425$$
$$426 \quad f(x) = ax \quad (11)$$
$$427$$

428

429 The Mean Squared Error (MSE) loss function is used in this study to calculate the difference
430 between the actual and the prediction value. This function computes the average squared
431 difference between the actual and predicted values using Equation (12) [43]. The MSE is also
432 used as the metric parameter for keeping track of performance measures (i.e., the objective
433 function of the NN).

$$434$$
$$435$$
$$436 \quad MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (12)$$
$$437$$

438

439 y_i = the actual value, \hat{y}_i = the predicted value, n = number of predicted values

440

441 The "learning rate" is a crucial hyperparameter in NN that determines the step size for
442 weight updates during the optimization process. Overfitting is a common issue when data is
443 similar to each other in time, and to avoid it, the "shuffle" parameter should be used while
444 training the model. The Adam optimizer is used for updating weights in the network, and it can
445 adjust the learning rate with the "decay" parameter. The network's weights are updated iteratively
446 using batches, with each batch representative of the dataset, and the batch size should be large
447 enough to include non-zero values. The dataset is divided into three splits of training, testing, and
448 validation, with each containing a different percentage of the data. The optimal values for the
449 hyperparameters require evaluating different values on the dataset.

450 The process of updating the weights of a neural network is iterative and occurs over many
451 epochs. An epoch consists of one forward and backward pass of the entire dataset, which is
452 usually too large to be fed into the network at once, so it is divided into batches. It is important to
453 choose a batch size that is representative of the dataset to prevent errors. The batch size
454 parameter should be large enough to include non-zero values in a batch since multiple output
455 values are often zero. The correct values for the batch size and epoch parameters should be
456 chosen by assessing different data values. Finally, the dataset is split into three sets: training,
457 testing, and validation, with percentages of 56%, 30%, and 14%, respectively.

458

459 **4. Case Study**

460 A case study was done using open-source trip and weather datasets to evaluate the proposed
461 framework. These datasets are reviewed in detail in the following sections.

462

463 **4.1. Trip Data**

464 This study uses open-source data provided by the Taxi and Limousine Commission (TLC)
465 available on the NYC Open Data website [46]. The dataset consists of 8.81 million New York

466 Green Taxi 2018 trip records. Each record includes pickup date and time, drop-off date and time,
467 trip's origin and destination zones ID, and other fields shown in Table 1 (redundant fields for
468 modeling such as tax and distance are ditched).

469 According to the City Zones dataset available on the NYC Open Data, zone IDs denoted in
470 Table 1 represent specific taxi zones defined by the Department of City Planning [47]. As shown
471 in Figure 2, there are 265 zones numbered from 1 to 265. Table 2 shows samples of this dataset.
472 This study uses this definition of zones for the case study instead of TAZs for specifying origins
473 and destinations.

474

475 **4.2. Weather Data**

476 Section 2 discussed that the weather conditions substantially impact daily travel behaviors.
477 This study uses an open-source weather dataset from the National Climatic Data Center (NCDC)
478 [48]. This dataset contains hourly recorded weather information, including date and time,
479 temperature, wind speed, wind direction, 1-hour liquid precipitation, 6-hour liquid precipitation,
480 and snow depth, as shown in Table 3.

481

482 **4.3. Data Preparation**

483 Now that the datasets are explained, procedures done on the datasets described in section 3
484 are summarized here. After setting constraints for each parameter, values exceeding these
485 constraints are omitted. Then, examining the remaining data reveals that there are apparently no
486 outliers remaining in the dataset, as they were possibly removed in the last step. However, there
487 would still be outliers while calculating average travel times, which will be removed based on
488 the Z-score. Moreover, additional constraints should be set for some parameters. For example,
489 calculated travel times with minimal positive values (e.g., less than a minute) may not be omitted
490 based on the initial constraint of being positive and also may not be detected as outliers. After
491 cleaning the datasets from possible errors, only the data for the first quarter of the year (90 days)
492 is used in this study to prevent inundating the network while training. So, indices of data rows
493 considering the available 265 zones are calculated as $265*265*90*24 = 151,686,000$ indices.

494

495 **4.4. Data Verification**

496 It is necessary to verify the validity of the acquired data to prevent possible errors before
497 training the network. Data verification can be done by examining the conformity of trip patterns
498 with previous studies. As a result, average trip patterns are investigated in this section to verify
499 the validity of the dataset. To investigate the trip patterns, average hourly trip counts were
500 computed for all weekends and weekdays in February 2018. To show the contrast between the
501 trip patterns during the holidays and non-holidays, hourly trip counts of 1st January 2018 (New
502 Year's Day) are computed as a sample. These patterns are plotted in Figure 3.

503 As illustrated in Figure 3, there are two peak hours in the morning and the evening for trip
504 counts in the weekday average trip pattern, which are 8 a.m. and 6 p.m., respectively. The
505 weekend average trip pattern shows that the morning peak hour is vanished (since there are no
506 work-based trips in the morning) and the midnight trip is increased significantly compared to the
507 average weekdays. The trip patterns from previous studies can be compared to similar hourly
508 patterns of taxi trips on weekdays and weekends to validate the results [30,49]. The main
509 difference between the holidays and non-holiday trip patterns (including non-holiday weekends)
510 is that there are numerous holidays, and generally, traffic patterns are changed based on the

511 holiday and related celebrations or rituals of the day. The reasons mentioned above substantially
512 impact trip patterns (e.g., travel destinations change vastly). For instance, the hourly trip pattern
513 on 1st January 2020 plotted in Figure 3 shows substantial differences compared to weekday and
514 weekend trip patterns.

515 Trip patterns between OD zones can also be inspected to see differences in the patterns
516 between weekends, workdays, and holidays. The weekday average trip pattern shows that trips
517 between zones 74 and 75, East Harlem North and East Harlem South neighborhoods, have the
518 most frequency at different day hours, including the morning and evening peak hours. According
519 to the Office of the New York State Comptroller report, East Harlem is mainly a residential
520 neighborhood with concentrated small businesses [50]. Weekends average trip pattern exposes
521 that trips between zones 41 and 42 (Central Harlem and Central Harlem North neighborhoods)
522 and internal trips of zone 7 (Astoria zone of borough Queens) have the highest frequency at
523 different hours of the day. Inspecting New York City's Zoning and Land Use Map [51] indicates
524 that Central Harlem and Astoria are commercial neighborhoods, including numerous recreational
525 places, specifically the Astoria. Internal trips of the Astoria neighborhood also showed the
526 highest frequency at different hours in the selected holiday. According to the City Zones dataset,
527 these paths are illustrated in Figure 4.

528 These results gave us good insights into the differences between weekdays, weekends, and
529 holiday trip patterns and the need to use binary parameters to address these variances. It can also
530 be derived that the demand in the origin and destination zones is a function of land use.
531 Therefore, it can be verified that the input data have rational patterns and can be used for training
532 the network.

533

534 **5. Results**

535 **5.1 Network Results**

536 “Keras is a deep learning API written in Python, running on top of the machine learning
537 platform TensorFlow. It was developed with a focus on enabling fast experimentation.” [52] This
538 package provides the required functions for training the network in Python. After evaluating
539 different values for the network parameters, a summary of the chosen values for the network's
540 parameters is given in Table 4. After manually inspecting the network's prediction accuracy,
541 these values are chosen by evaluating different values for each parameter.

542 As shown in Figure 5, it is observed that the error value converges to a relatively constant
543 value after performing several epochs. Consequently, the number of epochs is chosen to be 10.
544 The NN training results indicate the presence of $MSE = 0.5798$ after ten epochs, which is
545 reasonable. It should be noted that these results cannot be compared directly to results from the
546 trip distribution models, including the Gravity Model and the Fratar model, since these models
547 predict aggregated trips for a period of time. On the contrary, this paper proposed a method to
548 predict hourly trip counts. However, the hourly prediction may cause an increase in error, which
549 comes from numerous possible scenarios in each hour. The output results of each epoch can be
550 seen in Figure 5, and the loss reduction trend in Figure 6. The middle oscillations in Figure 6
551 indicate the beginning of a new epoch. The error reduction trend verifies that the network is
552 working correctly. As shown in Figure 5, loss values decrease rapidly at first and then slowly
553 after the third epoch.

554 It can be inferred from the results that the loss value is minimized after a limited number of
555 iterations, and there is no need to increase the number of epochs. The network performance is
556 then investigated on one million random samples from the test dataset. It is worth mentioning

557 that the predicted results are rescaled to the original values and rounded to the nearest integer
 558 since they represent trip counts. Test results are given in Table 5.

559 Test results show that the number of predicted zero values in the OD matrix and the total
 560 predicted trips perfectly match the actual values, and the Mean Squared Error of 0.0348 confirms
 561 this. R Squared value of 0.453 shows the model's acceptable fit, but possible reasons for the R
 562 Squared's relatively small value are discussed here. One million samples are approximately
 563 equivalent to 14 hours of trips since there are 70,225 possible paths (a path is a possible route
 564 between OD pair) between zones in an hour, and the results show that there are 15,100 trips in
 565 one million random samples of the trips. It can be deduced that there is an average of 1,100 trips
 566 for the available 70,225 paths in one hour, which means the average hourly trip count for each
 567 path is a small value. It can also be derived from the results that the average value for non-zero
 568 trip counts is approximately equivalent to 1.7 trips. Hence, slight deviations from the actual
 569 value can be due to the rounded predicted values (e.g., the predicted value of 3 for the actual
 570 value of 2), decreasing the Coefficient of Determination (R Squared) vastly. As a result, R
 571 Squared's small value can not necessarily represent the model's inadequate goodness of fit, and
 572 the MSE is a better quantifier to evaluate the model's goodness of fit. Suggested solutions to
 573 reduce the existing errors are given in the discussion.

574
 575 **5.2 Validating Network Results**

576 In this section, the NN's results are compared to the Gravity model to validate the results. As
 577 mentioned before, the NN predicted hourly trips, and the Gravity model generates aggregated
 578 trip predictions; therefore, these results cannot be compared directly. So, the results of the
 579 Gravity model should be compared to the aggregated results of the NN. Although this study aims
 580 to predict hourly flows, comparing the aggregated form of the results with the traditional models
 581 is compulsory for validation. The Gravity model's general form can be expressed in Equation
 582 (13) [53,54].

583
 584
 585
$$T_{ij} = \frac{P_i A_j F_i K_j \times f(c_{ij})}{\sum_v A_v F_i K_v \times f(c_{ij})} \quad (13)$$

586
 587
 588 T_{ij} = total trips between zones i and j

589 P_i = total trips produced by zone i

590 A_j = total number of trips attracted to zone j

591 v = set of 265 zones

592 $f(c_{ij})$ = decreasing function of the travel cost c_{ij}

593 F_i = balancing factor ensuring $\sum_j T_{ij} = P_i$

594 (14) K_j = balancing factor ensuring $\sum_i T_{ij} = A_j$

596 (15)

597

598 The travel cost function in Equation (13) (friction function) is any decreasing function of the
 599 travel cost (which is assumed to be the travel time in this study). Hence, the friction function can
 600 be considered as the power function shown in Equation (16) [55].

601
 602

$$603 \quad f(c_{ij}) = \frac{1}{(c_{ij})^n} \quad (16)$$

604
 605

606 c_{ij} = average travel time between zones i and j
 607 n = power variable
 608

609 Since the Gravity model requires total productions and attractions for the prediction period
 610 (i.e., the NN test data), they should be estimated using the trip generation models. The trip
 611 generation models require access to socio-economic data, which are assumed to be unavailable in
 612 this study. As a result, zones' productions and attractions are estimated by applying linear
 613 regression to the train data to predict the number of attracted (A_j) and produced (P_i) trips for the
 614 test set. The training and test data in this section cannot be the same as before since the Gravity
 615 model predicts aggregated trips for the prediction period. Consequently, hourly records of the
 616 dataset are aggregated into daily records. The test data includes 27 daily productions and
 617 attractions records for all zones (30% of the whole period), and the training data includes 63
 618 daily records. Although the trip generation estimation should be for the desired 27 days, the
 619 training and the test data are aggregated into data points of 9 days to increase accuracy. In other
 620 words, the test data is aggregated into three data points (each one including aggregated
 621 productions and attractions of all zones for nine days period), and the training data is aggregated
 622 into seven accumulated data points.

623 Then, linear regression is applied to each zone's seven data points to predict future
 624 productions and attractions. It is worth mentioning that the linear regression equation is
 625 calibrated for the productions and attractions of each zone separately. Hence, $265 \times 2 = 530$
 626 linear regression equations are calibrated, predicting three data points for future periods. Zones'
 627 production and attraction values are then compared to the actual values. Table 6 shows R
 628 Squared values of predictions for the zones' productions and attractions.

629 Linear regression results show a reasonable fit of the predicted productions and attractions
 630 with an R Squared of 0.99 and a reasonable error in predicting the total trips. The three predicted
 631 data points for each zone's attractions and productions are then aggregated to calibrate the
 632 Gravity model. As shown in Equation (14), the Gravity model requires average travel times for
 633 the forecasting period. The required travel times are calculated by averaging the non-zero travel
 634 times after removing the outliers, as described in section 3.3. It should be noted that the average
 635 travel times are calculated using the data from the 63-day training dataset. Since the zero average
 636 travel times cannot be used in Equation (16), the zone pairs with an average travel time of zero
 637 are assumed to have no trip interchanges. However, this assumption may increase the accuracy
 638 of the results since zone pairs with no trip interchanges in the 63 days training period are forced
 639 to have no trips in the future.

640 After implying the travel times in Equation (13), the F_i and K_j balancing factors are
 641 calibrated through an iterative process to ensure that the conditions expressed in Equation (14)
 642 and Equation (15) are met. The stop condition of this iterative process is as follows:
 643

$$644 \max \left\{ \max_{i \in v} \left(\left| 1 - \frac{P_i}{\sum_j T_{ij}} \right| \right), \max_{j \in v} \left(\left| 1 - \frac{A_j}{\sum_i T_{ij}} \right| \right) \right\} < 0.05 \quad (17)$$

645
 646 T_{ij} = total trips between zones i and j
 647 P_i = total trips produced by zone i
 648 A_j = total number of trips attracted to zone j
 649 v = set of 265 zones
 650

651 Then, the calibration process is done with different power variables, n , in Equation (16) to
 652 minimize the error. Table 7 shows the Gravity model results with different values for the power
 653 variable, including the iterations needed to meet the convergence condition expressed in
 654 Equation (17).

655 As shown in Table 7, $n = 2$ had the lowest error in prediction with an MSE of 2340 and an R
 656 Squared of 0.84. However, the aggregated NN results for the same period showed an MSE of
 657 less than 25. One point worth mentioning from the above table is that the R Squared value for the
 658 power variable of 6 is negative. While the R Squared name suggests that it may always range
 659 from 0 to 1, some exceptions may also be negative. In cases where the model predictions are not
 660 being compared to the observation that were used for calibrating the model, the Total Sum of
 661 Squared Errors component (SS_{res}) is not included in the Total Sum of Squares (SS_{tot}) [56].
 662
 663

$$664 R^2 = 1 - \frac{SS_{res}}{SS_{tot}} \quad (18)$$

665
 666
 667 SS_{res} = Total Sum of Squares of Residuals
 668 SS_{tot} = Total Sum of Squares
 669

670 Hence, as Equation (18) suggests, the R Squared value could also be negative in such cases.
 671 It should be noted that R Squared has been criticized for the lack of reliability as a measure of
 672 predictive accuracy [57]. Therefore, a more suitable measure of accuracy should be used to
 673 compare the prediction results, which is MSE in this case. Moreover, the NN uses MSE as the
 674 metric to optimize the learning process, and as a result, the network tries to reduce MSE in each
 675 iteration. The results showed that the NN had a clear superiority to the Gravity model in this
 676 case, although the purpose of this study was to predict hourly OD flows, and the Gravity model
 677 is unable to predict the flows on an hourly basis. Besides, forecasting the future trip distribution
 678 using the Gravity model required additional steps to estimate the future trip generations (i.e., the
 679 second step of the FSM), while the NN can predict trip distributions more precisely without
 680 requiring further steps. Another point worth mentioning is that, as mentioned before, the output

681 results of the gravity may not be directly compared to those of the NN. The reason is that the
682 Gravity model predicts aggregated trips for a period of time, while the NN in this study aims to
683 predict hourly trips. Although it can be argued that the aggregated trips of the NN output results
684 can be compared to the Gravity model results, given that the NN is optimized to predict the
685 hourly hours, this study is not focusing on such comparisons using visual descriptions (e.g.,
686 comparison plots)

687 688 **6. Conclusion and Discussion**

689 OD matrix prediction for a specific transit mode using traditional methods has always faced
690 numerous problems, including data collection. Using traditional methods requires data collection
691 for the trip generation and trip distribution steps of the FSM. This data includes users' socio-
692 economic characteristics and travel expenses information, which requires time-consuming and
693 costly collection methods, such as filling out questionnaires. Finally, due to the nature of these
694 collecting data methods, there is a significant error in the collected data, and it is also challenging
695 to update them periodically.

696 This paper aimed to facilitate this procedure using the data from data-driven transportation
697 systems. Prediction results showed proper fit and the logical dependence of the output data on
698 the input data. Other advantages of predicting trips using the NN compared to the traditional
699 modeling methods are considering more scenarios (weekends/holidays and more), quickly
700 updating the network with recent changes, and adequately forecasting OD flows on an hourly
701 basis.

702 It can be inferred from the results that there are considerable differences in the number of
703 trips between zones. As a result, some output values, which are numerically significant and
704 essential to be included in predictions, are detected as outliers and have insignificant impacts on
705 model training. As a potential research extension, paths between zones can be classified based on
706 their traffic volume (e.g., low traffic, medium traffic, and high traffic) and then modeled for each
707 category separately. Creating dummy variables indicating each category can also be done
708 instead.

709 It should also be noted that the results represent the predicted part of the demand that taxi
710 drivers could handle. In other words, there would be other trip demands exceeding the taxi
711 service supply; therefore, as there are no data for the unanswered demands in the dataset, the
712 calibrated model disregards such demands. A potential research extension includes datasets
713 containing users' requests to consider the drivers' unhandled trip requests, especially during peak
714 hours.

715 As described in section 4.4, holiday trips showed various patterns depending on the
716 occasion. A potential research direction is to model each type of trip discussed in this article
717 separately (e.g., separate modeling for weekdays and weekends) to increase accuracy. Using
718 algorithms to detect abnormal trip patterns (e.g., gatherings, special occasions that are not
719 officially registered, and social events) and separating them from other data used for training the
720 network can also improve results. Trip data used in this study included origin and destination
721 zones for each trip, including precise longitude and latitude of origins and destinations, resulting
722 in more accurate travel times and improved results.

723

724 **Conflict of Interest**

725 All the authors have no conflict of interest with the funding entity and any organization
726 mentioned in this article in the past three years that may have influenced the conduct of this
727 research and the findings.

728

729 **Acknowledgments and Declarations**

730 This research did not receive any specific grant from funding agencies in the public,
731 commercial, or not-for-profit sectors.

732

733 **References**

734

- 735 1. Mathew, T. V and Sharma, S., “Capacity expansion problem for large urban transportation
736 networks”, *J. Transp. Eng.*, **135**(7) (2009).
- 737 2. Halvorsen, A., Koutsopoulos, H. N., Ma, Z., and Zhao, J., “Demand management of
738 congested public transport systems: a conceptual framework and application using smart
739 card data”, *Transportation*, **47**(5) (2020).
- 740 3. McNally, M. G., “The Four-Step Model”, In *Handbook of Transport Modelling*, D. A.
741 Hensher and K. J. Button, Eds., Emerald Group Publishing Limited, **1**, pp. 35–53 (January
742 1, 2007).
- 743 4. Bhat, C. R., “Random utility-based discrete choice models for travel demand analysis”,
744 *Transp. Syst. Plan. Methods Appl.*, **10**(1), pp. 1–30 (2003).
- 745 5. Manski, C. F., “The structure of random utility models”, *Theory Decis.*, **8**(3) (1977).
- 746 6. Walker, J. and Ben-Akiva, M., “Generalized random utility model”, *Math. Soc. Sci.*, **43**(3),
747 pp. 303–343 (2002).
- 748 7. Liu, C., Susilo, Y. O., and Karlström, A., “Weather variability and travel behaviour – what
749 we know and what we do not know”, *Transp. Rev.*, **37**(6), pp. 715–741 (2017).
- 750 8. Rudloff, C., Leodolter, M., Bauer, D., Auer, R., Brög, W., and Kehnscherper, K.,
751 “Influence of Weather on Transport Demand: Case Study from the Vienna, Austria,
752 Region”, *Transp. Res. Rec.*, **2482**(1), pp. 110–116 (2015).
- 753 9. Tsapakis, I., Cheng, T., and Bolbol, A., “Impact of weather conditions on macroscopic
754 urban travel times”, *J. Transp. Geogr.*, **28**, pp. 204–211 (2013).
- 755 10. Litman, T., “Land Use Impacts on Transport How Land Use Factors Affect Travel
756 Behavior” (2008).
- 757 11. McNally, M. G., *The Four-Step Model*, Emerald Group Publishing Limited (2007).
- 758 12. Bouchard, R. J. and Pyers, C. E., “Use of Gravity Model for Describing Urban Travel”,
759 *Highw. Res. Rec.*, **88** (1965).
- 760 13. Long, G. D., *An Evaluation of the Gravity Model Trip Distribution*, Texas Transportation
761 Institute (1968).
- 762 14. McFadden, D. L., “Conditional Logit Analysis of Qualitative Choice Behavior”, *Front.*
763 *Econom.* (1974).
- 764 15. Louviere, J., Street, D., Carson, R., Ainslie, A., Deshazo, J. R., Cameron, T., Hensher, D.,
765 Kohn, R., and Marley, T., “Dissecting the Random Component of Utility”, *Mark. Lett.*,
766 **13**(3), pp. 177–193 (2002).

- 767 16. Ranganathan, P., Pramesh, C. S., and Aggarwal, R., “Common pitfalls in statistical
768 analysis: Intention-to-treat versus per-protocol analysis”, *Perspect. Clin. Res.*, **7**(3), pp.
769 144–146 (2016).
- 770 17. Cramer, J. S., “The Origins of Logistic Regression”, *SSRN Electron. J.* (2005).
- 771 18. Lindner, A., Pitombo, C. S., and Cunha, A. L., “Estimating motorized travel mode choice
772 using classifiers: An application for high-dimensional multicollinear data”, *Travel Behav.*
773 *Soc.*, **6**, pp. 100–109 (2017).
- 774 19. Duan, Z., Zhang, K., Chen, Z., Liu, Z., Tang, L., Yang, Y., and Ni, Y., “Prediction of city-
775 scale dynamic taxi origin-destination flows using a hybrid deep neural network combined
776 with travel time”, *IEEE Access*, **7**, pp. 127816–127832 (2019).
- 777 20. Zhang, J., Che, H., Chen, F., Ma, W., and He, Z., “Short-term origin-destination demand
778 prediction in urban rail transit systems: A channel-wise attentive split-convolutional neural
779 network method”, *Transp. Res. Part C Emerg. Technol.*, **124**, p. 102928 (2021).
- 780 21. Chu, K. F., Lam, A. Y. S., and Li, V. O. K., “Deep Multi-Scale Convolutional LSTM
781 Network for Travel Demand and Origin-Destination Predictions”, *IEEE Trans. Intell.*
782 *Transp. Syst.*, **21**(8), pp. 3219–3232 (2020).
- 783 22. Krishnakumari, P., van Lint, H., Djukic, T., and Cats, O., “A data driven method for OD
784 matrix estimation”, *Transp. Res. Part C Emerg. Technol.*, **113**, pp. 38–56 (2020).
- 785 23. Ke, J., Qin, X., Yang, H., Zheng, Z., Zhu, Z., and Ye, J., “Predicting origin-destination ride-
786 sourcing demand with a spatio-temporal encoder-decoder residual multi-graph
787 convolutional network”, *Transp. Res. Part C Emerg. Technol.*, **122**, p. 102858 (2021).
- 788 24. Lee, D., Derrible, S., and Pereira, F. C., “Comparison of Four Types of Artificial Neural
789 Network and a Multinomial Logit Model for Travel Mode Choice Modeling”, *Transp. Res.*
790 *Rec.*, **2672**(49) (2018).
- 791 25. Golshani, N., Shabanpour, R., Mahmoudifard, S. M., Derrible, S., and Mohammadian, A.,
792 “Modeling travel mode and timing decisions: Comparison of artificial neural networks and
793 copula-based joint model”, *Travel Behav. Soc.*, **10** (2018).
- 794 26. Brathwaite, T., Vij, A., and Walker, J. L., “Machine Learning Meets Microeconomics: The
795 Case of Decision Trees and Discrete Choice” (2017).
- 796 27. Xiong, X., Ozbay, K., Jin, L., and Feng, C., “Dynamic Origin–Destination Matrix
797 Prediction with Line Graph Neural Networks and Kalman Filter”, *Transp. Res. Rec.*,
798 **2674**(8), pp. 491–503 (2020).
- 799 28. Yaldi, G., Taylor, M. a. P., and Yue, W. L., “Using Artificial Neural Network in Passenger
800 Trip Distribution Modelling (A Case Study in Padang, Indonesia)”, *Proc. East. Asia Soc.*
801 *Transp. Stud.*, **2009**, pp. 105–105 (2009).
- 802 29. Hyland, M., Frei, C., Frei, A., and Mahmassani, H. S., “Riders on the storm: Exploring
803 weather and seasonality effects on commute mode choice in Chicago”, *Travel Behav. Soc.*,
804 **13**, pp. 44–60 (2018).
- 805 30. Dong, X., Wang, L., and Hu, B., “Analysis of spatio-temporal distribution characteristics of
806 passenger travel behaviour based on online ride-sharing trajectory data”, *J. Phys. Conf.*
807 *Ser.*, **1187**(5), p. 052055 (2019).
- 808 31. Regehr, J. D., Montufar, J., and Hernandez-Vega, H., “Traffic Pattern Groups Based on
809 Hourly Traffic Variations in Urban Areas”, *J. Transp. Inst. Transp. Eng.*, **7**(1), pp. 1–16
810 (2015).

- 811 32. Fujita, M., Yamada, S., and Murakami, S., “Time Coefficient Estimation for Hourly Origin-
812 Destination Demand from Observed Link Flow Based on Semidynamic Traffic
813 Assignment”, *J. Adv. Transp.*, **2017** (2017).
- 814 33. Venkatanarayana, R., Smith, B. L., and Demetsky, M. J., “Quantum-Frequency Algorithm
815 for Automated Identification of Traffic Patterns”, *Transp. Res. Rec.*, **2024**(1), pp. 8–17
816 (2007).
- 817 34. Yu, Q., Xie, Y., Li, W., Zhang, H., Liu, X., Shang, W.-L., Chen, J., Yang, D., and Yan, J.,
818 “GPS data in urban bicycle-sharing: Dynamic electric fence planning with assessment of
819 resource-saving and potential energy consumption increasement”, *Appl. Energy*, **322**, p.
820 119533 (2022).
- 821 35. Shang, W.-L., Chen, J., Bi, H., Sui, Y., Chen, Y., and Yu, H., “Impacts of COVID-19
822 pandemic on user behaviors and environmental benefits of bike sharing: A big-data
823 analysis”, *Appl. Energy*, **285**, p. 116429 (2021).
- 824 36. Shang, W.-L., Chen, Y., and Ochieng, W. Y., “Resilience Analysis of Transport Networks
825 by Combining Variable Message Signs With Agent-Based Day-to-Day Dynamic Learning”,
826 *IEEE Access*, **8**, pp. 104458–104468 (2020).
- 827 37. Shang, W., Han, K., Ochieng, W., and Angeloudis, P., “Agent-based day-to-day traffic
828 network model with information percolation”, *Transp. Transp. Sci.*, **13**(1), pp. 38–66
829 (2017).
- 830 38. Deri, J. A., Franchetti, F., and Moura, J. M. F., “Big data computation of taxi movement in
831 New York City”, *2016 IEEE Int. Conf. Big Data Big Data*, pp. 2616–2625 (2016).
- 832 39. Freire, J., Bessa, A., Chirigati, F., Vo, H., and Zhao, K., “Exploring What not to Clean in
833 Urban Data: A Study Using New York City Taxi Trips”, *IEEE Data Eng. Bull.*, **39**(2), pp.
834 63–77 (2016).
- 835 40. Patel, U., *NYC Taxi Trip and Fare Data Analytics Using BigData* (2015).
- 836 41. Martínez, L. M., Viegas, J. M., and Silva, E. A., “A traffic analysis zone definition: a new
837 methodology and algorithm”, *Transportation*, **36**(5), pp. 581–599 (2009).
- 838 42. Rousseeuw, P. J. and Hubert, M., “Robust statistics for outlier detection”, *WIREs Data Min.*
839 *Knowl. Discov.*, **1**(1), pp. 73–79 (2011).
- 840 43. Shanker, M., Hu, M. Y., and Hung, M. S., “Effect of data standardization on neural network
841 training”, *Omega*, **24**(4), pp. 385–397 (1996).
- 842 44. LeCun, Y., Bengio, Y., and Hinton, G., “Deep learning”, *Nature*, **521**(7553), pp. 436–444
843 (2015).
- 844 45. Nwankpa, C., Ijomah, W., Gachagan, A., and Marshall, S., “Activation Functions:
845 Comparison of trends in Practice and Research for Deep Learning” (2018).
- 846 46. Taxi, N. Y. and (TLC), L. C., “New York City Green Taxi Trip Data”, Available:
847 <https://www.nyc.gov/site/tlc/about/tlc-trip-record-data.page> (2018).
- 848 47. Department of City Planning’s Neighborhood Tabulation Areas (NTAs), N. Y. C. D.,
849 “NYC Taxi Zones”, Available: [https://catalog.data.gov/dataset/nyc-taxi-](https://catalog.data.gov/dataset/nyc-taxi-zones/resource/79e158e7-d158-4ace-b9d3-41ede473c76c)
850 [zones/resource/79e158e7-d158-4ace-b9d3-41ede473c76c](https://catalog.data.gov/dataset/nyc-taxi-zones/resource/79e158e7-d158-4ace-b9d3-41ede473c76c) (2019).
- 851 48. (NCDC), N. C. D. C., “NYC 2018 Hourly Surface Data”, Available:
852 <https://www.ncei.noaa.gov/products/land-based-station/integrated-surface-database> (2018).
- 853 49. Tang, J., Liu, F., Wang, Y., and Wang, H., “Uncovering urban human mobility from large
854 scale taxi GPS data”, *Phys. Stat. Mech. Its Appl.*, **438**, pp. 140–153 (2015).

855 50. Office of the New York State Comptroller, O., “An Economic Snapshot of the East Harlem
856 Neighborhood”, Available: [https://www.osc.state.ny.us/files/reports/osdc/pdf/report-9-](https://www.osc.state.ny.us/files/reports/osdc/pdf/report-9-2018.pdf)
857 2018.pdf (2018).

858 51. Office of City Planning, N. Y. C. D., “New York City’s Zoning & Land Use Map”, p. 5,
859 Available: <https://zola.planning.nyc.gov/> (2020).

860 52. Chollet, F., “Keras”, Available: <https://github.com/fchollet/keras> (2015).

861 53. Wilson, A. G., “Advances and problems in distribution modelling”, *Transp. Res.*, **4**(1), pp.
862 1–18 (1970).

863 54. Duffus, L. N., Sule Alfa, A., and Soliman, A. H., “The reliability of using the gravity model
864 for forecasting trip distribution”, *Transportation*, **14**(3), pp. 175–192 (1987).

865 55. Celik, H. M., “Sample size needed for calibrating trip distribution and behavior of the
866 gravity model”, *J. Transp. Geogr.*, **18**(1), pp. 183–190 (2010).

867 56. Willmott, C. J., “On the Validation of Models”, *Phys. Geogr.*, **2**(2), pp. 184–194 (1981).

868 57. Schemper, M., “Predictive accuracy and explained variation”, *Stat. Med.*, **22**(14), pp. 2299–
869 2308 (2003).

870

871 **Captions**

872 Figure 1. Structure of the Neural Network in this study

873 Figure 2. New York City taxi zones

874 Figure 3. Weekday, Weekend and Holiday hourly trip patterns. Holiday trip data as of 1st January 2018, Weekend
875 and Weekday trip data are average trip counts of all weekends and all weekdays of February 2018 respectively.

876 Figure 4. Paths with highest trip counts

877 Figure 5. Optimum loss values in each epoch

878 Figure 6. Loss values in each iteration

879 Table 1 - Trip record fields sample

880 Table 2 – New York City zones dataset sample

881 Table 3 - Weather record sample

882 Table 4. Network parameters values

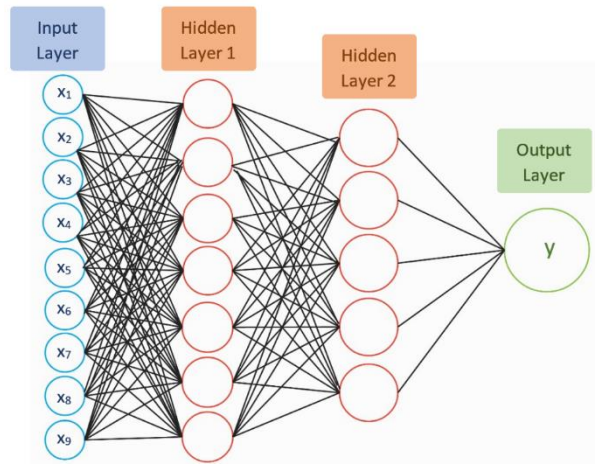
883 Table 5. Test results on one million random samples

884 Table 6 - Trip generation linear regression results

885 Table 7 - The Gravity model results

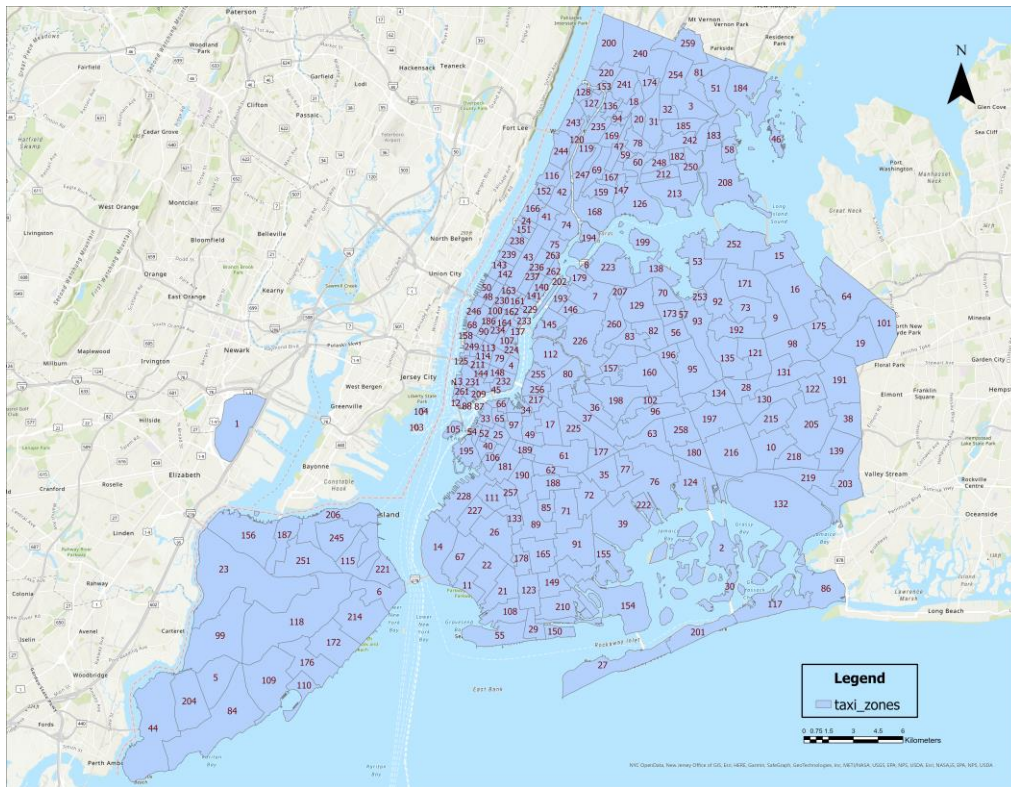
886

887



889
890

Figure 2. Structure of the Neural Network in this study



891
892
893

Figure 2. New York City taxi zones

894

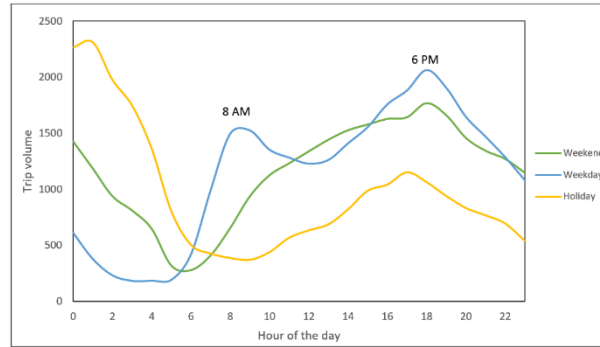


Figure 3. Weekday, Weekend and Holiday hourly trip patterns
 Holiday trip data as of 1st January 2018, Weekend and Weekday trip data are average trip counts of all weekends and all weekdays of February 2018 respectively.

895
896

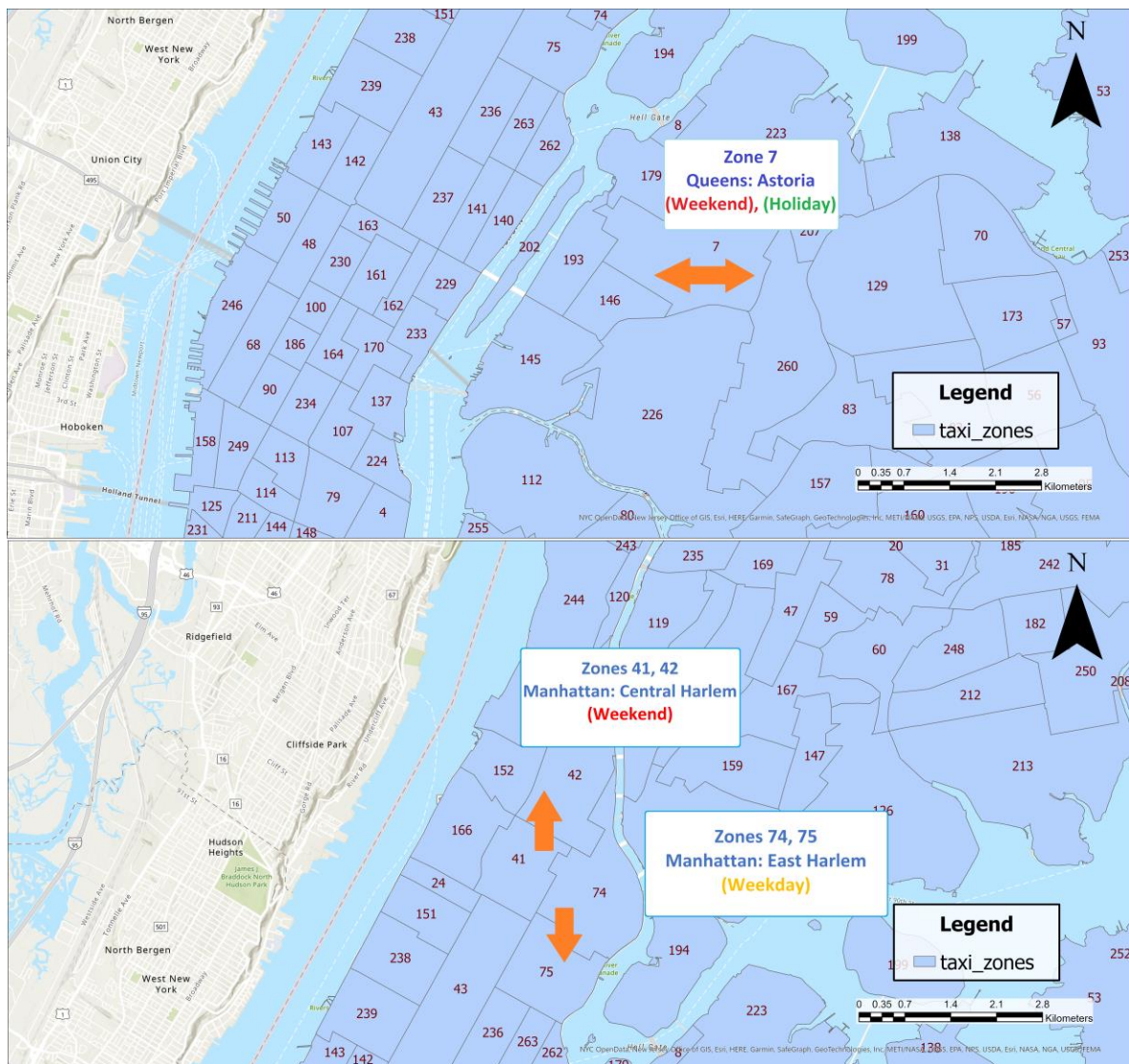


Figure 4. Paths with highest trip counts

897
898
899

900
901
902

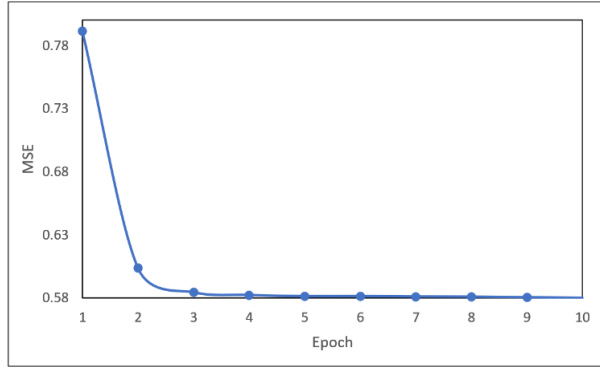


Figure 5. Optimum loss values in each epoch

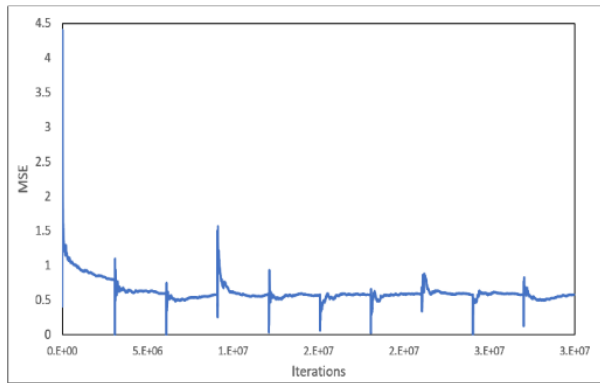


Figure 6. Loss values in each iteration

903
904
905
906
907

Table 8 - Trip record fields sample

Field	Pickup Date Time	Dropoff Date Time	PU Location ID	DO Location ID	Passenger Count	Trip Distance
Description	Passenger pick-up date and time	Passenger drop-off date and time	Destination zone ID	Origin zone ID	Number of passengers	Distance traveled from origin to destination
Sample	01/01/2018 12:18:50 AM	01/01/2018 12:24:39 AM	43	75	2	3.5 mi

908
909
910

Table 9 – zones dataset

Zone ID	12	46	94	165
Location	Manhattan-Battery Park	Manhattan-Chinatown	The Bronx-Fordham South	Brooklyn-Midwood

New York City sample

911
912
913
914

915

Table 10 - Weather record sample

Field	YR-- MODAHRMN	DIR	SPD	SKC	TEMP	PCP01	SD
Description	Year-Month- Day-Hour- Minute in GMT	Wind direction in compass degrees	Wind speed (mph)	Sky cover (Nominal classification)	Temperature in Fahrenheit	1-hour liquid precipitation in inches	Snow depth in inches
Sample	201801281151	990	6	SCT	49	0.01	0.0

916

917

Table 11. Network parameters values

Parameter	Learning rate	Decay	Batch size	Number of epochs
Value	0.0005	$1 \cdot 10^{-6}$	256	10

918

919

Table 12. Test results on one million random samples

	Total Trips	Number of zero values	R Squared	MSE
Predicted Values	14721	991384	0.453	0.0348
Actual Values	15127	991324		

920

921

922

Table 13 - Trip generation linear regression results

	Total Trips		R Squared
	Predicted	Actual	
Zones Productions	692585	693364	0.9921
Zones Attractions	692565	693364	0.9922

923

924

925

Table 14 - The Gravity model results

Power Variable (n)	Number of zero values		R Squared	MSE	Iterations
	Actual	Predicted			
	1	19112			
2	17780		0.8439	2340.261	11
3	16584		0.7714	3248.935	13
4	13502		0.4755	7866.803	16
5	10164		0.1162	13254.1566	19
6	7950		-0.2248	18370.005	22

926

927