

Predictive Distance-based Road Pricing — Designing Tolling Zones through Unsupervised Learning

Antonios F. Lentzakis^{a,*}, Ravi Seshadri^b, Moshe Ben-Akiva^c

^a*NCS Hub, NCS Group, 5 Ang Mo Kio Street 62, Singapore 569141*

^b*DTU Management, Technical University of Denmark, Bygningstorvet, 116, 116A, Kongens Lyngby, Denmark*

^c*Intelligent Transportation Systems Lab, MIT, 77 Massachusetts Avenue, Room 1-181, Cambridge, M.A., 02139, U.S.A.*

Abstract

Congestion pricing is a standard approach to mitigate traffic congestion in a number of urban networks around the world. The advancement of satellite technology has spurred interest in distance-based congestion pricing schemes, which obviate the need for fixed infrastructure such as gantries that are used in area- and cordon-based pricing. Moreover, distance-based pricing has the potential to more effectively manage traffic congestion. In the context of distance-based congestion pricing, we propose the use of sparse subspace clustering methods employing Elastic Net optimization (SCELE) and Orthogonal Matching Pursuit (SSCOMP), as well as two hierarchical density-based clustering methods, (OPTICS, HDBSCAN*) for the derivation of tolling zones. These tolling zone derivations are then used within a simulation-based framework for real-time predictive distance-based toll optimization to examine network congestion and performance of the tolling schemes. Within this framework, for a given derivation of tolling zones, tolling function parameters are optimized in real-time using a simulation-based Dynamic Traffic Assignment (DTA) model. Guidance information generation is integrated into the predictive optimization framework and behavioral responses to the information and tolls along dimensions of departure time, route, mode, and trip cancellation are explicitly modeled. For the evaluation of network performance we make use of Travel Speed Index (TSI) data from the real-world Boston Central Business District urban network and demonstrate that tolling zones derived from the sparse subspace cluster-

*Corresponding author: aflentz@mit.edu

ing are effective means of operationalizing real-time distance-based toll optimization schemes and can positively impact overall network performance, showing improvements in average travel time and social welfare relative to the baseline.

Keywords: Distance-based Toll Optimization, Sparse Subspace Clustering, Density-based Clustering

1. Introduction

Traffic congestion is a serious issue world-wide, which results in large costs to travelers, the environment and economy. Congestion was estimated to result in a total of 5.5 billion hours of time delay and 2.9 billion gallons of fuel expenditure in urban areas in the United States between 2000 and 2010 (Litman, 2019) and the costs of congestion were projected to increase from \$121 billion in 2011 to \$199 billion in 2020. Mitigating congestion is always a high-priority and also impacts transportation network reliability, driver’s comfort and traffic safety. Congestion pricing is a standard approach for congestion mitigation that influences traveler behavior along several dimensions: trip making and frequency, mode, destination, time of day, route, and so on. Traditional approaches to congestion pricing include facility-based and area-based schemes (de Palma and Lindsey, 2011) that rely on physical infrastructure such as gantries or gates for vehicle detection. Unfortunately, the reliance on fixed physical infrastructure makes it difficult to modify or relocate the charging areas or zones. Moreover, these schemes can result in inefficiency in terms of congestion mitigation since they do not differentiate toll charges based on the associated externalities or congestion caused due to differing distances traveled or time spent in congestion. The aforementioned disadvantages of area- and facility-based pricing and the advancement of Global Navigation Satellite Systems (GNSS) have focused attention on usage-based tolling wherein toll charges depend on the distance-traveled or the time spent in congestion (see Smith et al. (1994) and Bonsall and Palmer (1997) for a detailed discussion on the comparative performance of distance- and time-based schemes). Singapore is in the process of transitioning to such a GNSS-based electronic road-pricing scheme (ERP2) (LTA, 2016, 2021). Distance-based schemes may be operationalized by dividing the urban area into zones and charging a distance-based toll such that the tariff varies across zones and by time-of-day. The motivation for the use of tolling zones (instead of a single distance-based scheme over the entire network) is that it provides the flexibility to adjust the tolling rates based on road-type and congestion levels, thereby improving overall efficiency gains.

Past research on area and cordon-based real-time toll optimization has typically applied reactive approaches (where the optimization of tolls is not based on forecasts of future traffic conditions, but

rather on prevailing traffic conditions) for small corridor networks and there are few studies that adopt a predictive approach in the context of large networks (Gupta et al., 2016, 2020). A more detailed discussion of cordon and area-based real-time toll optimization may be found in Gupta et al. (2020). As noted previously, in contrast with cordon- and area-based schemes, distance-based tolling schemes involve partitioning the network into zones, and levying a toll within each zone that is a function of distance traveled (linear toll functions are considered in Gu et al. (2018); Yang et al. (2012); Zhu and Ukkusuri (2015), and piece-wise linear functions are used in Liu et al. (2014); Meng et al. (2012); Sun et al. (2016)). Distance-based toll optimization problems have largely been formulated as simulation-based optimization problems (Gu and Saberi, 2019b; Gu et al., 2018; Lentzakis et al., 2020), non-linear programs (Yang et al., 2012) and mathematical programs with equilibrium constraints or MPEC (Liu et al., 2014; Meng et al., 2012), which are solved by global optimization approaches (Liu et al., 2014), meta heuristics (Lentzakis et al., 2020; Meng et al., 2012), reinforcement learning (Zhu and Ukkusuri, 2015) and feedback controllers (Gu and Saberi, 2019b; Gu et al., 2018). With the exception of Lentzakis et al. (2020), these approaches are based on prevailing network conditions (i.e., they are reactive as opposed to proactive), and do not consider elastic demand or the integration of guidance information generation.

Several studies have also examined the partitioning of networks utilizing flow, speed and density data (Ji and Geroliminis, 2012; Lentzakis et al., 2014; Saeedmanesh and Geroliminis, 2017) for the design of traffic management schemes utilizing the Network Fundamental Diagram (NFD) concept. Although area- and cordon-based pricing has been studied in great detail (Geroliminis and Levinson, 2009; Simoni et al., 2015; Zheng et al., 2016, 2012), distance-based pricing in particular has only recently received attention on idealized networks (Daganzo and Lehe, 2015), using nested regions (Gu et al., 2018) and at the link-level (Simoni et al., 2019). With the exception of Lentzakis et al. (2020), there has been limited research on systematic approaches for the derivation of tolling zones within distance-based toll optimization strategies. Due to the increasing significance of distance-based road pricing in traffic network management and operations, this paper addresses the problem of how to define tolling zones and proposes the application of sparse subspace clustering methods to define parsimonious sets

54 of tolling zones. The performance of these methods is evaluated within a framework for real-time
55 toll optimization which generates predictive optimized distance-based toll strategies combined with
56 guidance information. This paper contributes to the existing literature in the following respects:

- 57 1. We apply sparse subspace and hierarchical density-based clustering methods for the derivation of
58 tolling zones that utilize location coordinates and travel speed indices (TSI) as features. The key
59 advantage of using sparse subspace clustering techniques is that they enable the effective use of
60 high-dimensional temporal network performance data (for example, travel speeds at a resolution of
61 five minutes) directly in the clustering algorithm. This provides a potentially promising alternative
62 to the procedure proposed in [Lentzakis et al. \(2020\)](#) where the clustering algorithm is applied to
63 a single aggregate measure of network performance (over the entire peak period) for each link. In
64 this paper, we focus specifically at the performance of the different clustering algorithms and the
65 implications for toll design/policy.
- 66 2. The proposed clustering methods are evaluated using a framework for real-time distance-based
67 predictive toll optimization on the Boston CBD network and yield insights into their performance
68 and suitability for deployment wherein one of our primary goals is minimization of computational
69 effort.

70 **2. Framework for Predictive Distance-based Toll Optimization**

71 In this section, we summarize the real-time distance-based predictive toll optimization framework
72 (more details may be found in [Lentzakis et al. \(2020\)](#)), the optimization problem formulation, the
73 proposed clustering methods for tolling zone derivation and the algorithmic solution for the optimization
74 problem.

75 *2.1. Framework*

76 The framework, shown in Figure 1, uses DynaMIT2.0 - a simulation-based Dynamic Traffic Assign-
77 ment (DTA) system developed at the MIT Intelligent Transportation Systems Lab ([Ben-Akiva et al.](#),

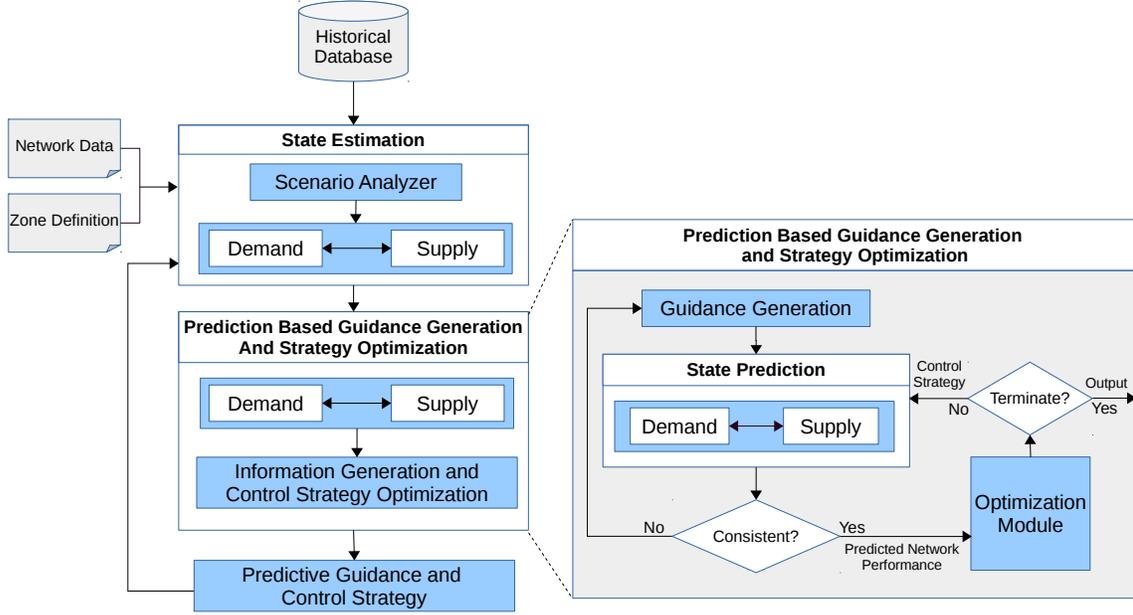


Figure 1: Real-time distance-based predictive toll optimization framework

78 2010; Lu et al., 2015b). DynaMIT2.0 employs a rolling horizon approach involving two key modules,
 79 state estimation and state prediction. The state estimation process uses a combination of historical data,
 80 real-time traffic surveillance data, and prevailing network control strategies (such as distance-based toll
 81 optimization) to estimate the current state of the network. It used detailed models of demand (pre-trip
 82 models of departure time, route and mode choice), supply (mesoscopic traffic simulator that com-
 83 bines speed-density relationships and a deterministic queuing model) and their interactions. Following
 84 this, the state prediction module generates forecasts of traffic conditions for a pre-specified prediction
 85 horizon (origin-destination demands and supply parameters are forecasted for the future using an au-
 86 toregressive process). The strategy optimization and guidance generation modules in conjunction use
 87 the state predictions to first, optimize control strategies for the prediction horizon and second, generate
 88 guidance information (traveler information) for the prediction horizon. The evaluation of candidate
 89 control strategies makes use of network predictions and guidance information that are consistent, i.e.,

90 the guidance information is as close as possible to actual predicted network travel times (see Figure 1
 91 and Ben-Akiva et al. (2010) for more on this aspect of consistency).

92 3. Problem Formulation and Solution

93 In this section, we describe the optimization problem formulation (based on the framework described
 94 in Section 2) including details of the demand model within the DTA system, and the solution approach.

95 3.1. Context and Tolling Function Definition

96 We represent the transportation network of interest as a directed graph $\mathcal{G} = (\mathcal{N}, \mathcal{A})$, where \mathcal{N}
 97 denotes the set of n network nodes and \mathcal{A} denotes the set of m links. The network is partitioned into
 98 $l = 1 \dots L$ tolling zones, where every zone l is defined by a subset of network links $\mathcal{A}_l \subseteq \mathcal{A}$. For each
 99 zone l , we define a tolling function $\phi_l(\boldsymbol{\theta}_l^t, D_l)$ that maps distance traveled within the zone l, D_l to the
 100 toll amount; $\boldsymbol{\theta}_l^t$ is a vector of parameters that defines the tolling function in time interval t . Further, it
 101 is assumed that the toll payable in a zone is bounded, i.e $\tau_{LB} \leq \phi_l(\boldsymbol{\theta}_l^t, D_l) \leq \tau_{UB}, \forall l = 1, 2, \dots, L \forall t =$
 102 $1, 2, \dots, T$.

103 Denote the length of the state estimation interval in DynaMIT2.0 by Δ (usually 5 minutes) and
 104 assume that the prediction horizon is composed of H such intervals so that the size of the prediction
 105 horizon is $H\Delta$. We assume that the prediction horizon and the optimization horizon are identical.
 106 Further, the tolling function parameters do not vary within a given time interval of size Δ and these
 107 tolling intervals coincide with DynaMIT2.0 estimation intervals. For an arbitrary estimation interval
 108 $[t_0 - \Delta, t_0]$, let $\boldsymbol{\theta}^h = (\boldsymbol{\theta}_1^h, \boldsymbol{\theta}_2^h \dots \boldsymbol{\theta}_L^h)$ represent the vector of tolling function parameters for the time
 109 period $[t_0 + (h - 1)\Delta, t_0 + h\Delta]$ where $h = 1, \dots, H$. Accordingly, for the current optimization horizon,
 110 the decision variables are $\boldsymbol{\theta} = (\boldsymbol{\theta}^1, \boldsymbol{\theta}^2, \dots, \boldsymbol{\theta}^H)$.

111 Implementing a system with complex zone-based tariffs that vary every five minutes is likely to
 112 impose unreasonable burdens on drivers that may compromise acceptability of the system. An added
 113 issue is that drivers may not have a viable alternative if for example they suddenly find themselves
 114 entering a zone where the tariff has increased substantially. Hence, we assume that drivers are charged

115 the predicted toll that the system provides to them at the point of departure (more precisely, the point
 116 at which they make their decision, which may be up to 15 minutes prior to their actual departure).
 117 The underlying premise – justifiable given our rolling horizon design – is that the predictions of the
 118 toll in the future do not deviate appreciably from the actual implemented tolls. The rolling horizon
 119 framework is demonstrated in Figure 2

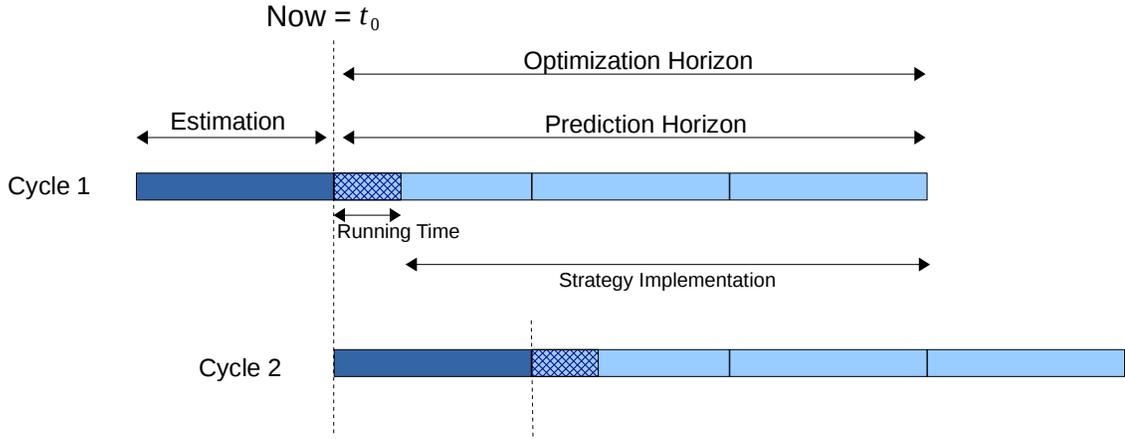


Figure 2: Rolling Horizon Approach for Tolling Function Optimization

120 Consider the set of vehicles $v = 1, \dots, V$ that are on the network during the prediction horizon
 121 $[t_0, t_0 + H\Delta]$. For each vehicle v , we denote the experienced trip travel on its chosen route by tt^v
 122 and the predictive guidance information by $\mathbf{tt}^g = (\mathbf{tt}_i^g; \forall i \in A)$, where \mathbf{tt}_i^g represents a vector of time
 123 dependent travel times for link i . Note that the vehicle travel times $\mathbf{tt} = (tt^v; v = 1, \dots, V)$ are obtained
 124 from the state prediction module of DynaMIT2.0, which we characterize through a single constraint
 125 that represents the coupled demand and supply simulators as:

126

$$G(\mathbf{x}^P, \gamma^P, \mathbf{tt}^g, \theta) = \mathbf{tt} \quad (1)$$

127

128 Where \mathbf{x}^P, γ^P represent the forecasted demand and supply parameters for the prediction horizon, and

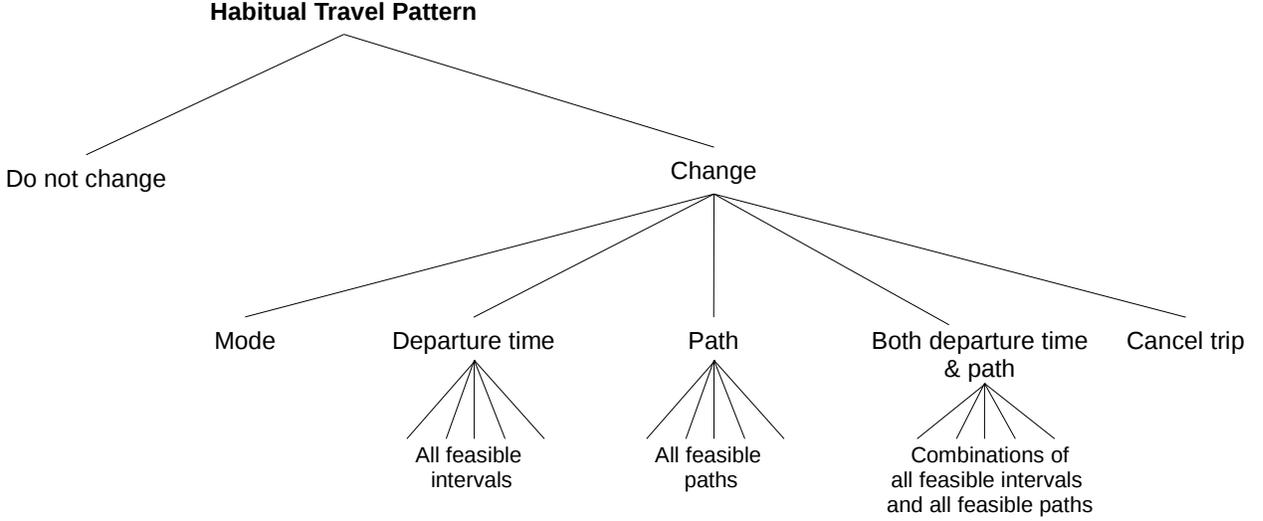


Figure 3: Pre-trip Behavior Model

129 θ is the vector of tolling function parameters. As noted previously, the state prediction module ensures
 130 consistency between \mathbf{tt}^g and \mathbf{tt} .

131 3.2. Pre-trip Behavioral Model with Elastic Demand

132 The pre-trip response of users to the travel time guidance and distance-based tolls is modeled using
 133 a path-size nested logit model with heterogeneous value of time (illustrated in Figure 3) that captures
 134 decisions of mode choice, trip cancellation, departure time and path (notation is provided in Table
 135 1). We provide a brief description of the model here, for completeness (more details may be found in
 136 [Lentzakis et al. \(2020\)](#)).

137 In response to pre-trip information and tolls, a traveler may alter his/her habitual travel pattern,
 138 which may include changing mode, canceling trip, changing departure time or path, or changing depar-
 139 ture time and path. This results in elastic total demand w.r.t. traffic congestion. The options of mode
 140 modeled are private car (drive alone) and public transit. The utility of change to transit for vehicle v
 141 is given by:

Table 1: Pre-trip Model - Abbreviations

Abbreviation	Variable
β_{CM}	Alternative Specific Constant (ASC) for change of mode to transit
β_{CT}	ASC for canceling trip
β_{CDT_d}	ASC for departure in time interval d
c_m^v	monetary cost for traveling with non-private (transit) mode
c_{dp}^v	toll charge for departure via path p in interval d
c_p^v	toll for switching to path p
t_m^v	travel time associated with non-private (transit) mode
t_{dp}^g	travel time (guidance) for departure via path p in interval d
t_p^g	travel time (guidance) for switching to path p
$at_{d'p}^{hab}$	arrival time (habitual)
at_{dp}^g	arrival time (predicted) for departure via path p in time interval d
β_c^v	monetary cost coefficient
β_t^v	travel time coefficient
β_E	schedule delay early coefficient
β_L	schedule delay late coefficient
PS_p	path size variable
C_*	utility relating to number of left turns/signalized intersections and path length
ε_*	error component

142

$$U^v(CM) = \beta_{CM} + \beta_c^v c_m^v + \beta_t^v t_m^v + \varepsilon_m \quad (2)$$

143

144 The utility of departing at time interval d and choosing path p for vehicle v is given by:

$$\begin{aligned} U_{dp}^v &= \beta_{CDT_{dp}} + \beta_c^v c_{dp}^v + \beta_t^v t_{dp}^g \\ &+ \beta_E \max(at_{d'p'}^{hab} - at_{dp}^g, 0) + \beta_L \max(at_{dp}^g - at_{d'p'}^{hab}, 0) \\ &+ \log(PS_p) + C_{dp} + \varepsilon_{dp} \end{aligned} \quad (3)$$

where :

$$c_{dp}^v = \sum_{l=1}^L \phi_l(\theta_l^{t_{v,l}}, D_l^v),$$

145

146 and $t_{v,l}$, D_l^v denote the predicted time of entry of vehicle v into zone l and the total distance traveled
 147 by vehicle v in zone l , respectively. Note that if there are a total of N combinations of path and depar-
 148 ture time choices in the choice set, the alternative specific constant β_{CDT_d} can only appear in $(N - 1)$
 149 utilities. The utility of canceling trip altogether is given by:

150

$$U^v(CT) = \beta_{CT} + \varepsilon_{CT} \quad (4)$$

151

152 Thus, the probability of vehicle v choosing alternative c within the choice set C is given by:

153

$$P^v(c|C) = \frac{e^{\mu V_c^v}}{\sum_{a \in C} e^{\mu V_a^v}} \quad (5)$$

154

155 where V_c^v is the systematic utility given by $V_c^v = U_c^v - \varepsilon_c$ and μ is a scale parameter . The en-route
 156 choice model defines response of users in terms of path-choice to the toll and predictive travel time
 157 guidance. It is also formulated as a multinomial path size logit model where the utility of switching to

158 path p is given by:

159

$$U_p^v = \beta_c^v(tc_p^v) + \beta_t^v(tt_p^g) + \log(PS_p) + C_p + \varepsilon_p$$

where :

(6)

$$tc_p^v = \sum_{l=1}^L \phi_l(\theta_l^{t_{v,t}}, D_l^v)$$

160

161 Note that owing to the design of the distance-based tolls, which require that users are charged upfront
 162 at the beginning of a trip, we assume that en-route changes to the path are not made.

163 3.3. Optimization Formulation

164 The objective function for the toll optimization problem, formulated from the standpoint of the
 165 traffic regulator, is total social welfare (SW), which is the sum of the consumer surplus and the pro-
 166 ducer surplus. In this context, the consumer surplus (CS) is defined as the sum of the experienced
 167 utilities across all travelers, derived at the end of each simulation run, and the producer surplus is
 168 the net revenue, denoted by TP, which is simply the toll revenue minus variable costs (fixed costs are
 169 ignored), $TP = TR - VC$. We assume that the variable costs are a proportion of the toll revenue (the
 170 proportionality factor is denoted by $\alpha < 1$). Thus, the social welfare is given by:

$$\begin{aligned} SW &= CS + TP \\ &= CS + (TR - VC) \\ &= \sum_{v=1}^V \frac{U^v}{|\beta_c^v|} + \left[(1 - \alpha) \times \sum_{v=1}^V c^v \right] \end{aligned} \tag{7}$$

171

172 The absolute value of β_c^v is used to translate CS into dollar equivalents. The distance-based toll opti-
 173 mization problem is formulated as a simulation-based optimization problem in Equation 8, where the
 174 objective is social welfare, the decision variables are the vector of tolling function parameters for the
 175 current optimization horizon, and the constraints are toll bounds and the DTA model system.

176

$$\max_{\boldsymbol{\theta}} \left[\sum_{v=1}^V \frac{U^v}{|\beta_c^v|} + (1 - \alpha) \times \sum_{v=1}^V c^v \right]$$

(8)

$$\text{s.t.} \quad G(\mathbf{x}^{\mathbf{P}}, \gamma^{\mathbf{P}}, \mathbf{tt}^{\mathbf{g}}, \boldsymbol{\theta}) = \mathbf{tt}$$

$$\tau_{LB} \leq \phi_l(\boldsymbol{\theta}_l^h, D_l^v) \leq \tau_{UB}, \forall v = 1, 2, \dots, V; l = 1, 2, \dots, L; h = 1, 2, \dots, H$$

177 The upper and lower bounds τ_{LB}, τ_{LB} are imposed to allow for tolling function values within a safe
 178 and acceptable range, suitable for real-life implementations.

179 3.4. Solution Algorithm

180 Due to the highly non-convex nature of the objective function in 8, we apply a real-coded Genetic
 181 Algorithm (GA) to solve the optimization problem in 8. More details on the GA algorithm may be
 182 found in [Lentzakis et al. \(2020\)](#). Computational performance is enhanced by utilizing parallelization
 183 wherein the evaluations of different candidate solutions within an iteration of the GA are performed in
 184 parallel.

185 4. Tolling Zone Design through Unsupervised Learning

186 For most tolling-related implementation decisions, at least in the USA, tolling zone boundaries are
 187 subject to extreme political scrutiny, Environmental Justice reviews, exemptions for residents of certain
 188 areas, etc. Unsupervised machine learning utilizes historical data to reveal patterns, similarities or hid-
 189 den structure and can contribute to changing the current state-of-affairs, with regards to tolling zone
 190 design. In this paper, we posit that it would be beneficial for a city or highway operator to rely on un-
 191 supervised learning approach in any sort of real-world setting, for expeditious implementation. Even in
 192 the case that stakeholder involvement is mandatory, tolling zone derivation through unsupervised learn-
 193 ing can significantly augment the decision-making process. Unsupervised learning can be approached
 194 through different techniques such as clustering, association rules, and dimensionality reduction. Our

195 focus will be on clustering. One of the inputs with significant impact on our distance-based tolling sys-
196 tem performance is the tolling zone derivation. This input specifies which links belong to each tolling
197 zone and the number of zones. Each tolling function $\phi_l(\theta_l^h, D_l^y)$ corresponds to one tolling zone. Past
198 literature for partitioning urban traffic networks used datasets based on speed, flow, density (Gu and
199 Saberi (2019a); Ji and Geroliminis (2012); Lentzakis et al. (2014); Saeedmanesh and Geroliminis (2017))
200 and, more recently, marginal cost toll data (Lentzakis et al., 2020). In our case, the travel speed index
201 (TSI) is used, a widely used quantitative indicator that employs link speed normalization (Li and Xiao,
202 2020), given the fact that identical link speed levels might reflect different traffic conditions. Speed
203 information for toll setting is currently used in a similar fashion as in Singapore’s ERP system, (Lehe,
204 2019). It should be noted that the decision to use travel speed indices, rather than marginal cost tolls
205 (MCT) used in Lentzakis et al. (2020), has to do with the fact that, in this work, one of our main goals
206 was to reduce computational effort, both during data preprocessing and the predictive distance-based
207 toll optimization framework implementation, placing real-world applicability at the forefront. Should
208 circumstances allow it, the possibility of using MCT as a feature should definitely be explored.

209 4.1. Clustering Approaches

210 Elhamifar and Vidal (2009), inspired by compressed sensing (Lee et al., 2007), introduced Sparse
211 Subspace Clustering (SSC), which makes use of the self-expressiveness property to construct the affin-
212 ity matrix (which quantifies the extent of pairwise similarity between a set of data points). Self-
213 expressiveness (Elhamifar and Vidal, 2013) describes the fact that a data point found in a union of
214 subspaces can be represented as the linear combination of other data points. Based on the com-
215 puted affinity matrix, spectral clustering is applied to derive the underlying subspaces. While subspace
216 clustering methods have been used extensively for, among others, temporal video segmentation and
217 switched system identification (Bako, 2011; Rao et al., 2009), only recently, has this technique come to
218 the attention of the transportation research community. Zhang et al. (2019) employed SSC to classify
219 spatiotemporal taxi patterns with regards to their passenger searching behavior. For our experiments
220 we selected to compare two Sparse Subspace Clustering variants, SSCEL and SSCOMP, employing

221 Elastic Net optimization (You et al., 2016a) and Orthogonal Matching Pursuit (You et al., 2016b)
 222 respectively, against two well-known hierarchical density-based clustering methods, OPTICS (Ankerst
 223 et al., 1999), (Ordering Points To Identify the Clustering Structure), and HDBSCAN* (Campello et al.,
 224 2013), (Hierarchical Density-Based Spatial Clustering of Applications with Noise).

225 4.1.1. Sparse Subspace Clustering Methods

226 For the Sparse Subspace Clustering application, we selected two variants, SSCEL and SSCOMP. We
 227 exploited the property of self-representation to learn the affinity matrix, to be subsequently used in
 228 our implementation of spectral clustering. As noted previously, data self-expressiveness (Elhamifar and
 229 Vidal, 2013) describes the fact that a data point found in a union of subspaces can be represented as
 230 the linear combination of other data points, expressed through the following optimization problem:

$$\begin{aligned}
 & \min_{\mathbf{C}} \|\mathbf{C}\|_1 \\
 & \text{s.t.} \\
 & \mathbf{X} = \mathbf{X}\mathbf{C} \\
 & \text{diag}(\mathbf{C}) = 0
 \end{aligned} \tag{9}$$

231

232 Where $\mathbf{X} \in \mathbf{R}^{D \times N}$ is the data point matrix and $\mathbf{C} \in \mathbf{R}^{N \times N}$ is the self-expression coefficient matrix.
 233 In practice, however, solving N such problems over N variables may be computationally expensive for
 234 large N . Instead, the optimization problem is expressed as follows:

$$\begin{aligned}
 & \min_{\mathbf{c}_j} \|\mathbf{x}_j - \mathbf{X}\mathbf{c}_j\|_2^2 \\
 & \text{s.t.} \\
 & \|\mathbf{c}_j\|_0 \leq \mathbf{k} \\
 & \text{diag}(\mathbf{C}) = 0
 \end{aligned} \tag{10}$$

235 We can now efficiently solve the above problem using the Orthogonal Matching Pursuit algorithm, as
 236 described in You et al. (2016b). Orthogonal Matching Pursuit selects a single column of \mathbf{X} each time,

237 \mathbf{x}_j , such that the absolute value of the dot product with the residual \mathbf{c}_j is maximized and the coefficients
 238 are computed until k columns are selected. Subsequently, we learn the affinity matrix \mathbf{W} through data
 239 self-representation as: $\mathbf{W} = |\mathbf{C}| + |\mathbf{C}^T|$. Alternatively, we may employ Elastic Net regularization
 240 for scalable subspace clustering. Following [You et al. \(2016a\)](#), we used an active set algorithm that
 241 efficiently solves the elastic net regularization subproblem, which follows below, by capitalizing on the
 242 geometric structure of the elastic net solution:

$$\min_{\mathbf{c}_j} \lambda \|\mathbf{c}_j\|_1 + \frac{1-\lambda}{2} \|\mathbf{c}_j\|_2^2 + \frac{\gamma}{2} \|\mathbf{x}_j - \mathbf{X}\mathbf{c}_j\|_2^2 \quad (11)$$

243 Where $\lambda \in (0, 1]$ and $\gamma > 0$. In the majority of solution approaches for the Subspace Clustering
 244 problem, after learning the affinity matrix, spectral clustering is applied to the resulting matrix to
 245 derive the final clustering.

246 4.1.2. Hierarchical Density-based Clustering Methods

247 Hierarchical density-based clustering methods are gaining traction among the research community,
 248 exhibiting robustness during parameter selection and being able to cope with clusters characterized
 249 by large inter-cluster density variability, unlike their non-hierarchical predecessor, DBSCAN ([Schubert
 250 et al., 2017](#)).

251 OPTICS utilizes hyperparameters ϵ and κ , representing the maximum ball radius with each data point at
 252 its center and the minimum density threshold, respectively. Assuming a metric space (X, d) comprising
 253 of a set of data points $X = \{x_1, x_2, \dots, x_n\}$, a data point x is considered to be a core point with respect
 254 to ϵ and κ if its ϵ -neighborhood $N_\epsilon(x)$ contains a minimum of κ data points. Two core points x_i, x_j
 255 are ϵ -reachable with respect to ϵ and κ if they are both contained within each others ϵ -neighborhood.
 256 Two core points x_i, x_j are density-connected with respect to ϵ and κ if they are directly or transitively
 257 ϵ -reachable. A cluster is the largest possible group of data points, where each two points are considered
 258 connected in terms of density. In OPTICS data points are assigned a core distance $d_{\text{core}}^{\epsilon, \kappa}(x)$ to the κ -th
 259 nearest neighbor, for varying degrees of density. The reachability-distance $d_{\text{reach}}^{\epsilon, \kappa}(x_i, x_j)$ is the maximum
 260 between the core distance of x_i and the distance between data points x_i, x_j . A single global ϵ' value is

261 used to extract a flat clustering.

262 HDBSCAN* is similar to OPTICS with parameter $\epsilon = \infty$ and a different technique, based on cluster
263 stability, is utilized for flat clustering. In the case of HDBSCAN*, we have $d_{\text{core}}^{\kappa}(x_i)$ representing the
264 κ -th nearest neighbor distance. For a fixed κ and a range of possible ϵ values, the mutual reachability
265 distance $d_{\text{mreach}}^{\kappa}(x_i, x_j)$ is used to generate a complete hierarchy of clusterings. Thus, for any fixed ϵ
266 value, the clustering produced by DBSCAN at a given level in the hierarchy is the clustering obtained
267 for the corresponding ϵ value.

268 The selected hierarchical density-based clustering methods result in clusterings where some data
269 points are considered noise. A feasible derivation of tolling zones, however, must involve the assignment
270 of all data points to clusters. In order to address this issue, we perform a secondary assignment where
271 all noise data points are assigned to the closest clusters (using Euclidean distance).

272 4.2. Clustering Performance Metrics

273 As clustering performance metrics, the Silhouette Coefficient (SC) (Rousseeuw, 1987) and the
274 Davies-Bouldin index (DB) (Davies and Bouldin, 1979) were selected. SC is the average for the entire
275 dataset of the silhouette, which measures cohesion and separation for each cluster and ranges from
276 $[-1, 1]$, where -1 represents an inappropriate clustering (within-cluster variability is large and between
277 cluster variability is small), 0 represents overlapping clusters and 1 represents highly dense clustering.
278 DB is a function of the ratio of intra-cluster scatter to inter-cluster separation. DB values closer to 0
279 indicate a better clustering result.

280 4.3. Clustering Results

281 It would be preferable that the selected clustering methods produce clustering results that are of
282 high quality, according to our previously presented internal evaluation indices, but also do not preclude
283 any sort of practical application, due to high computational cost, incurred on the distance-based toll
284 optimization framework. This translates into a static tolling zone derivation (non-varying during the
285 simulation), with a reasonably low number of tolling zones. For our dataset, besides spatial coordinates,

286 we decided to use the travel speed index (TSI) for each link as an additional feature, calculated as follows:

$$\text{TSI}_i = 1 - \frac{\nu_i}{\nu_i^f} \quad (12)$$

287 Where ν_i, ν_i^f the link speed and free flow speed for link i respectively. Simulated speed data at the
288 segment and link level, obtained from a calibrated DynaMIT2.0 model of the Boston CBD (Lu et al.,
289 2015a), were used to derive the tolling zone derivations. In the case of static partitioning schemes
290 derived offline, a preferable alternative to using the average of TSI across specific intervals, as is the
291 case for our hierarchical density-based clustering approaches, would be to use the TSI values for all time
292 intervals, i.e., the entirety of our dataset, since self-representation, an integral part of sparse subspace
293 clustering, is amenable for use of datasets with spatiotemporal attributes (Hashemi and Vikalo, 2018;
294 Pham et al., 2012). The Boston CBD network, shown in Figure 4, has 846 nodes, 1746 links, 3085
295 segments, 5057 lanes and 13080 Origin-Destination pairs.

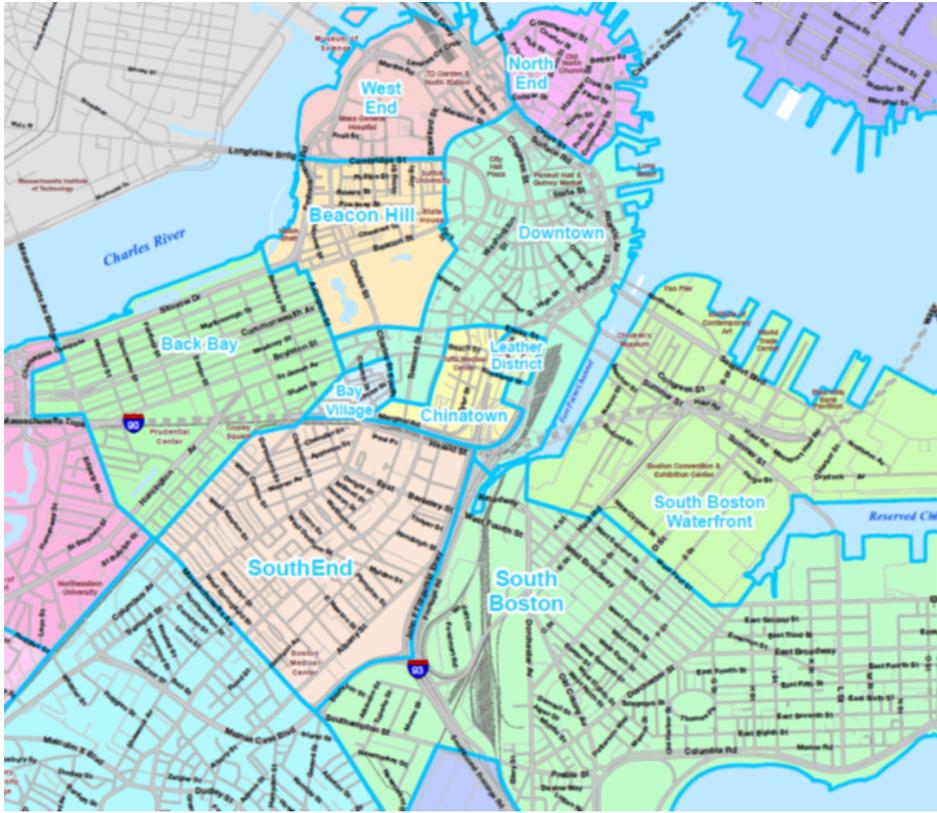


Figure 4: Boston CBD Network

SSCEL-TSI



(a) 2 zones derived using SSCEL for feature TSI, with $SC=0.114$, $DB=3.925$

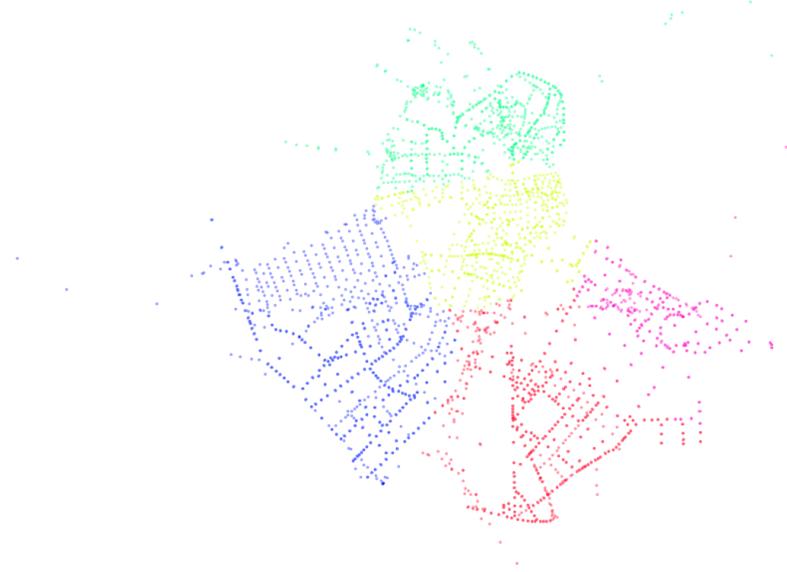
SSCOMP-TSI



(b) 2 zones derived using SSSCOMP for feature TSI, with $SC=0.198$, $DB=1.276$

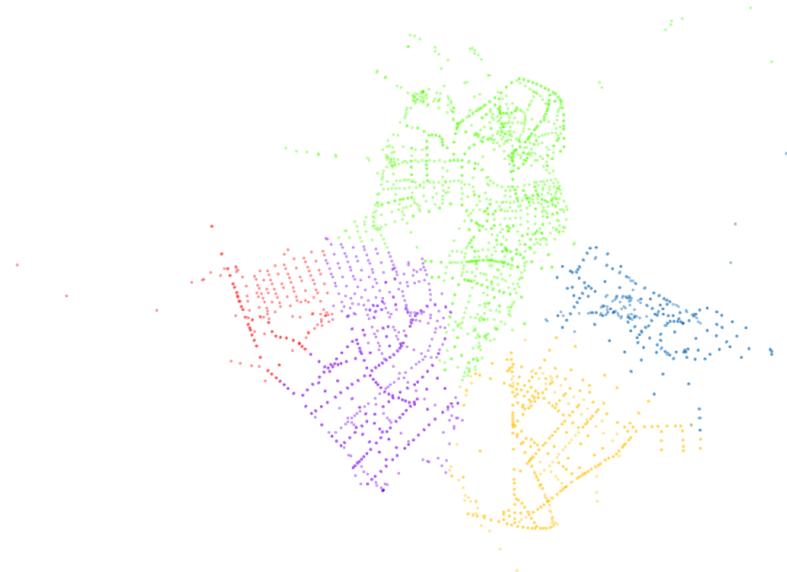
Figure 5: Clustering results and tolling zones (Sparse Subspace Clustering)

OPTICS-TSI



(a) 5 zones derived using OPTICS for feature TSI, with $SC=0.387$, $DB=0.858$

HDBSCAN-TSI



(b) 5 zones derived using HDBSCAN* for feature TSI, with $SC=0.400$, $DB=0.776$

Figure 6: Clustering results and Tolling zones (Hierarchical Density-based Clustering)

296 5. Experiments: Boston CBD Network

297 5.1. Experimental Design

298 In order to investigate the impact tolling zone derivations –which are derived from unsupervised
299 learning methods– have on the performance of adaptive distance-based congestion pricing schemes,
300 when applied on an urban network, experiments are conducted on the Boston CBD network illustrated
301 in Figure 4. A linear tolling function is considered with lower and upper bounds on the toll charged
302 in each zone (i.e $\phi_l(\boldsymbol{\theta}_l^t, D_l) = \theta_{l1}^t + \theta_{l2}^t D_l$; and $0 \leq \phi_l(\boldsymbol{\theta}_l^t, D_l) \leq 1.5$). The simulation period is from
303 06:00-09:00 covering the morning peak. As noted earlier, historical demand and supply parameters
304 are obtained from prior offline calibrations of DynaMIT2.0 for the Boston Central Business District
305 network (Azevedo et al., 2018). The estimation interval is 5 minutes and the prediction horizon is 30
306 minutes. The Boston network we consider contains 846 nodes, 1746 links, 3085 segments, 5057 lanes
307 and 13080 origin-destination pairs.

308 The performance measures are calculated for the population of vehicles with habitual departure
309 time within 06:00-09:00 (these drivers may later change the departure time in response to the traffic
310 conditions). A *warm-up* period of 15 min is used and the last 15 min of the simulation is a *cool-*
311 *down* period without toll optimization to ensure that all the vehicles with habitual departure time in
312 06:00-09:00 finish their trips.

313 The mean and standard deviation of the value of time are S\$23.5 and S\$5.75 respectively. The cost
314 coefficient for each vehicle is calculated from the lognormally distributed sampled value of time. The
315 parameters of the pre-trip choice model are summarized in Table 2.

316 A total of six scenarios are considered, which are summarized in Table 3. All scenarios involve
317 dynamic tolls computed using the framework described in Section 2, with the exception of the base
318 scenario (**B0**) which is the *No Toll* case. Recall that the state estimation interval is 5 minutes implying
319 that in the case of distance-based pricing schemes the tolling function parameters vary every 5 minutes.
320 The simulations were run using Ubuntu Linux on an HPC Cluster, with 5x60 cores and 5x250GB RAM.

321 The base scenario **B0** was calibrated to replicate prevailing traffic conditions in the Boston CBD

Table 2: Pre-trip Behavioral Model - Parameters

Parameter	Value
β_{CM}	-0.5
β_{CT}	-12
β_{CDT_1}	-0.12
β_{CDT_2}	-0.79
β_{CDT_3}	-1.15
β_{CDT_4}	- 1.65
β_t^v	-0.008
β_E	-0.004
β_L	-0.016

322 (refer to [Azevedo et al. \(2018\)](#) for more details). However, given the fact that a No Toll base scenario
 323 may not provide specific information regarding what portion of the performance uplift stems from these
 324 novel distance-based tolling schemes, rather than the inherent effects of using tolling to internalize the
 325 congestion externality, we are also considering a comparison scenario **B1**, (termed **UNIREG-TSI**),
 326 which employs predictive distance-based tolling on the the Boston CBD network as a unitary region.
 327 Scenarios **B2**, **B3**, **B4**, **B5** employ predictive distance-based tolling and differ only in the derivation
 328 of the tolling zones. In scenario **B2** (termed **SSCEL-TSI**), tolling zones are defined based on TSI
 329 data using SSCEL. In scenario **B3** (termed **SSCOMP-TSI**), tolling zones are defined based on TSI
 330 data using SSCP. In scenario **B4** (termed **OPTICS-TSI**), tolling zones are defined based on TSI
 331 data using OPTICS, and finally, in scenario **B5** (termed **HDBSCAN*-TSI**) they are based on TSI
 332 data using HDBSCAN*. Scenarios **B0-B5** are evaluated on three performance measures, total social
 333 welfare (SW), consumer surplus (CS) and average travel time (TT) to capture overall societal benefits,
 334 together with the impact on individual travelers.

Table 3: Simulation scenarios

Scenario	Tolling Scheme	Description
B0	No Toll	<i>No tolling scheme in place</i>
B1	Predictive distance-based	<i>Tolling zone encompassing entire network: UNIREG-TSI</i>
B2	Predictive distance-based	<i>Tolling zones derived from: SSCEL-TSI</i>
B3	Predictive distance-based	<i>Tolling zones derived from: SSCOMP-TSI</i>
B4	Predictive distance-based	<i>Tolling zones derived from: OPTICS-TSI</i>
B5	Predictive distance-based	<i>Tolling zones derived from: HDBSCAN*-TSI</i>

335 5.2. Results

336 The performance measures for all simulation scenarios are summarized in Table 4, the differences
337 in SW and CS (in \$ amounts) of scenarios **B1-B5** relative to the base scenario B0 are presented in
338 Figure 7a, and the relative performance in terms of average travel time (% improvement) over the base
339 scenario B0 is illustrated in Figure 7b. From Table 4, **B1-B5** exhibit an increase between \$182623.5
340 - \$206866.5 and \$64814.9 - \$127062.4, for SW and CS respectively, relative to **B0**. The average SW
341 gain per traveller, relative to the No Toll case is found to be around \$1.69 for those acquired via sparse
342 subspace clustering and around \$2.15 for tolling zone derivations acquired via hierarchical density-based
343 clustering.

344 Observe that all the scenarios yield a positive consumer surplus indicating that net user benefits
345 are positive even prior to any use of the toll revenues. This is a surprising finding and in contrast with
346 several past studies that have estimated negative user benefits (for example Eliasson and Mattsson
347 (2006) and De Palma et al. (2005)). We conjecture that this is a result of several factors. First, as
348 noted by Van Den Berg and Verhoef (2011)), in the case when there is heterogeneity in the value of time
349 (and values of schedule delay), the net user benefits may depend in large part on the extent and nature
350 of heterogeneity. In experiments on a variant of the standard bottleneck model including departure
351 time choice and a transit alternative (with heterogeneity in value of time, schedule delay, early and

late), [Chen \(2022\)](#) find that when the coefficient of variation in the value of time exceeds around 0.5, the net user benefits start to become positive (even before accounting for distribution of toll revenues). Second, our system integrates the provision of consistent guidance information with the optimization of tolls. These two factors coupled with the high levels of initial congestion may be the reason why we observe positive net user benefits even prior to a redistribution of toll revenues. Similar tests across different network topologies and spatio-temporal congestion patterns are required to determine whether this finding is a peculiarity of our context and network. The significant variation of differences in both welfare and CS across the five schemes confirms that the performance of distance-based tolling schemes is appreciably affected by the definition of the tolling zones. First, observe that the two sparse subspace clustering approaches yield quite different clusters and varying outcomes in terms of both CS and welfare. Scenario **B2** (SCCEL) yields the second lowest overall welfare, which is only marginally higher than scenario **B1** where the entire network is treated as a single zone. The reason for the relatively poor performance is two-fold. First, as is apparent in [Figure 5a](#), SCCEL results in clusters that lack spatial compactness. In other words, zones are 'non-contiguous' and links in different parts of the network belong to the same zone (links belonging to the red cluster or zone in particular). This clearly poses an issue in the toll optimization, since the toll design includes a two-part tariff where the fixed component is charged during each entry into a new zone (in other words each time a zone boundary is traversed). Overall, it results in the fixed part of the tariff being optimized at a much lower level than in the case when the zones are spatially compact (SSCOMP in Scenario **B3** and Scenarios **B4**, **B5**). Interestingly, the low tolls charged result in a high consumer surplus, comparable with the best performing scenario since the travel time gains and reductions in schedule delay costs are still significant.

The second reason for the poor performance of SCCEL in Scenario **B2** may be attributed to the clusters themselves. Observe that the key difference in the clusters or zones between Scenario **B2** and Scenario **B3** (which yields a significantly larger welfare) is that Scenario **B3** clearly demarcates the Back Bay region from the rest of Boston ([Figure 4](#)) whereas this is not the case in Scenario **B2**. The Back Bay region contains the Prudential center, which is a major attractor of trips in the morning peak and hence, arguably, the zone definitions in Scenario **B3** are more meaningful. This is also evident from

379 the clustering performance metrics which clearly indicate that the clusters are more homogeneous in
380 the case of Scenario **B3** than **B2** (SC of 0.198 versus 0.114).

381 Turning to the hierarchical density-based clustering approaches, we observe that both OPTICS
382 (Scenario B4) and HDBSCAN* (Scenario B5) yield meaningful clusters/zone definitions. Both dis-
383 tinguish the densely residential South Boston region (red cluster in Figure 6a and orange cluster in
384 Figure 6b; see also Figure 4) from the commercial South Boston Waterfront (pink cluster in Figure 6a
385 and blue cluster in Figure 6b). In Scenario B4, the commercial downtown region (dark green cluster
386 in Figure 6a) is separated from the more residential North End and West End regions (light green
387 cluster in Figure 6a). These three regions are all combined into a single zone in Scenario B5 (green
388 cluster in Figure 6b). The most notable difference in the clusters between B4 and B5 and one that
389 most likely leads to the significant performance difference is that Scenario **B5** clearly demarcates the
390 Back Bay region from the residential South End region (red and purple clusters in Figure 6b) unlike
391 Scenario B4. As discussed earlier, this appears to be the reason for Scenario B5 yielding the largest
392 gains in social welfare and consumer surplus. Scenario B5 also yields the clusters with links that are
393 internally homogeneous (SC of 0.4). Note that when using speed for clustering as we have done, the
394 more homogeneity within the clusters in an of itself does not appear to guarantee superior performance
395 in terms of welfare. This is evident when comparing Scenarios B3 and B4; B4 yields superior metrics
396 in terms of clustering performance but yields lower overall welfare. This underscores the importance
397 of checking the reasonableness of the clusters themselves using context specific knowledge of demand
398 patterns, land-use etc.

399 Notably, the TT performance improvement illustrated in Figure 7b, relative to the base case **B0** is
400 substantial in all schemes. It would appear that the Boston CBD area would benefit from an application
401 of a predictive distance-based tolling scheme, with average travel time TT improvements of up to 52%
402 (relative to **B0**). However, we do caution that the large travel time improvements may also be, in part,
403 due to a large number of short 'crossing' trips that are an artifact of modeling only the CBD area.

404 While it should be stated that while all the predictive distance-based tolling schemes yield substantial
405 network performance benefits when compared to the No Toll scenario, scenario **B5** with HDBSCAN*-

406 based tolling zone derivation yields the largest welfare gains. The observed welfare increase comes from
 407 the reduction in schedule delay costs and low travel times, which may be attributed to the efficient
 408 internalization of travel externality-associated costs through distance-based tolling. This in fact applies
 409 to all distance-based schemes considered in the experiments.

410 Scenario **B3** with only 2 tolling zones derived from SSCOMP resulted in comparable levels of
 411 performance to the Scenario **B5**, and in cases where computational effort poses a significant hurdle
 412 for practical implementation, it would be preferable to use SSCOMP. Although the HDBSCAN*-based
 413 tolling zone derivation leads to the best results, it is also more computationally intensive, due to the
 414 large number of tolling function parameters that require optimization. The overall computational time
 415 for scenarios B1-B3 were around 4 hours. Given we are simulating the 6-9 AM peak period, this does not
 416 yet achieve real-time performance for a 5-minute horizon. However, this could be attained by increasing
 417 the parallelization or marginally reducing the number of GA generations during the optimization. In
 418 case of scenarios B4, B5 where the number of zones are larger, the run times were significantly higher
 419 at 12 hours. In this case, to achieve real-time performance we would need to switch to a 15-minute roll
 420 period. For more details on computational considerations we refer the reader to [Gupta et al. \(2020\)](#).

Table 4: Performance measures

	Scenarios				
Metrics	<i>B1</i>	<i>B2</i>	<i>B3</i>	<i>B4</i>	<i>B5</i>
SW (\$)	181792.0	182623.5	205345.1	194744.3	206866.5
CS (\$)	64814.9	126194.4	110504.2	84083.7	127062.4
TT (s)	168.0	172.2	152.3	156.3	147.9

421 In Figure 8a, the Empirical Cumulative Distribution Function (ECDF) of the total toll charge
 422 (for the population of vehicles) for scenarios **B1, B2, B3, B4, B5** is presented. It is evident that
 423 the majority of the traveler population (almost 90%) pay total tolls no higher than \$3 for any of the
 424 distance-based pricing schemes. Further, the overall magnitude of toll charges in the case of scenario

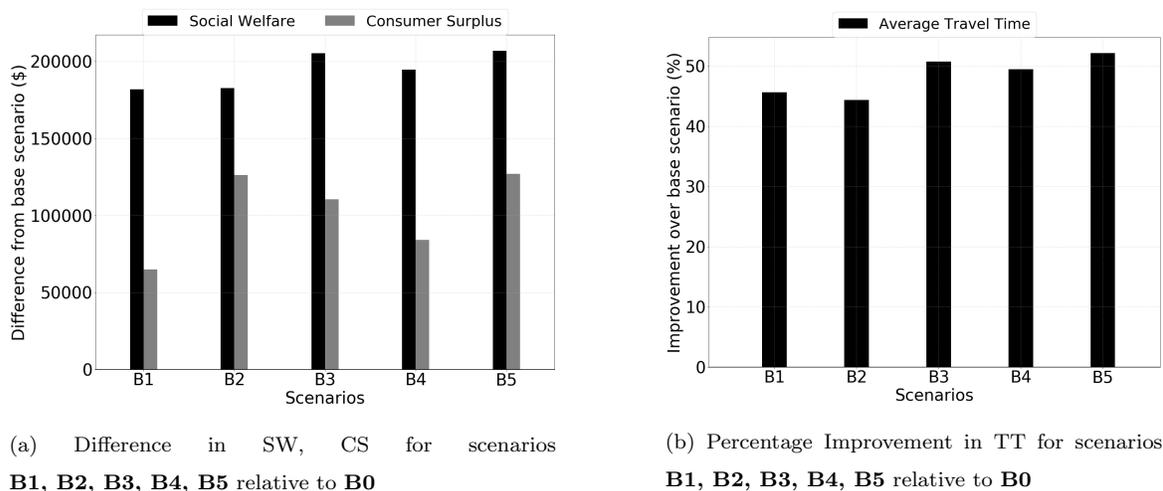
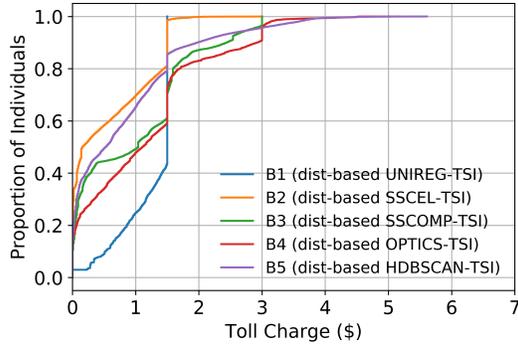


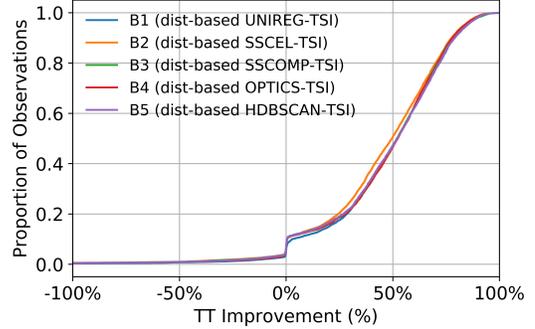
Figure 7: Performance results for scenarios **B1**, **B2**, **B3**, **B4**, **B5** relative to **B0**

425 **B4** is consistently higher than that of scenarios **B3**, **B5**, which happen to be the best performing
 426 scenarios. On the other hand, the overall magnitude of toll charges in scenario **B2** is consistently lower
 427 than that of scenarios **B3**, **B5**, which could be explained by the fact that the corresponding tolling zone
 428 derivation suffers in terms of spatial compactness, thus leading to lower tolling efficiency. For scenario
 429 **B1**, where the entire BCBD network is treated as a single tolling zone, we can observe that more than
 430 60% of the population is charged the toll upper bound (1.5\$), which leads to inequitable charging, but
 431 also higher revenue. The reason for this lies in the fact that, unlike the case of scenarios **B2-B5**, there
 432 is only one zone for the vehicles to traverse.

433 As is evident from Figure 8b, demonstrating the ECDF of travel time improvement per OD-pair, for
 434 less than 10% of the traveler population, travel times are equal or lower for scenario **B0**, as compared
 435 to scenarios **B1-B5**. Up to 90% of the traveler population benefits from lower travel times, in scenarios
 436 **B1-B5** employing distance-based tolling methods, compared to base scenario **B0**. It is also evident that
 437 the largest proportion of the traveler population subset that benefits from lower travel times corresponds
 438 to **B2**, though followed very closely by **B1**.



(a) ECDF of total Toll Charge values for scenarios **B1**, **B2**, **B3**, **B4**, **B5**



(b) ECDF of TT improvement per OD-pair over **B0** for scenarios **B1**, **B2**, **B3**, **B4**, **B5**

Figure 8: Empirical Cumulative Distribution Functions for scenarios **B1**, **B2**, **B3**, **B4**, **B5** relative to **B0**

439 *5.3. Iterating between toll zone definition and toll optimization*

440 Recall that the motivation behind using unsupervised learning for the toll zone definition is to
 441 decouple the problems of toll definition and toll value optimization. This avoids having to solve a
 442 complex mixed-integer programming problem for the design of the distance-based scheme. Second,
 443 and more importantly, the decoupling of the two problems also serves to provide a useful separation
 444 between what is performed *offline* and what is performed *online*. Specifically, the proposed framework
 445 involves setting the tolling zones offline and then optimizing the toll values in real-time every five or
 446 fifteen minutes (within for example, a traffic management system). In this context, it may be desirable
 447 to re-evaluate the zone definitions periodically, say every month or every quarter (as is done in the
 448 current ERP system in Singapore for the setting of the toll rates). In this setting, a loop from the toll
 449 optimization and the toll design would be beneficial.

450 In order to do so, we redo the clustering exercise in two ways. First, we compute an implied per-
 451 distance toll rate for each link and time interval from the optimized tolling function parameters obtained
 452 via the predictive distance-based toll optimization framework. We then use the resulting toll values
 453 from each scenario as a feature and perform the clustering once again with SSCEL,SSCOMP,OPTICS
 454 and HDBSCAN*, respectively. Note that clearly, since tolling function parameters are identical for

455 all links in the same zone by construction, this can only yield an aggregation of the original zone
456 definitions. Nevertheless, it serves to examine robustness of the original zone definitions. Second,
457 we redo the clustering using the travel speed indices obtained after the application of the optimized
458 tolls. Interestingly, in both cases we observe that the original clustering results and metrics are quite
459 robust (see Figures 9–12). However, in the case that they are not, this procedure could in principle be
460 performed iteratively and the zone definitions could be updated.

SSCEL-TOLL



(a) 2 zones derived using SSCEL for feature TSI, with $SC=0.118$, $DB=3.716$

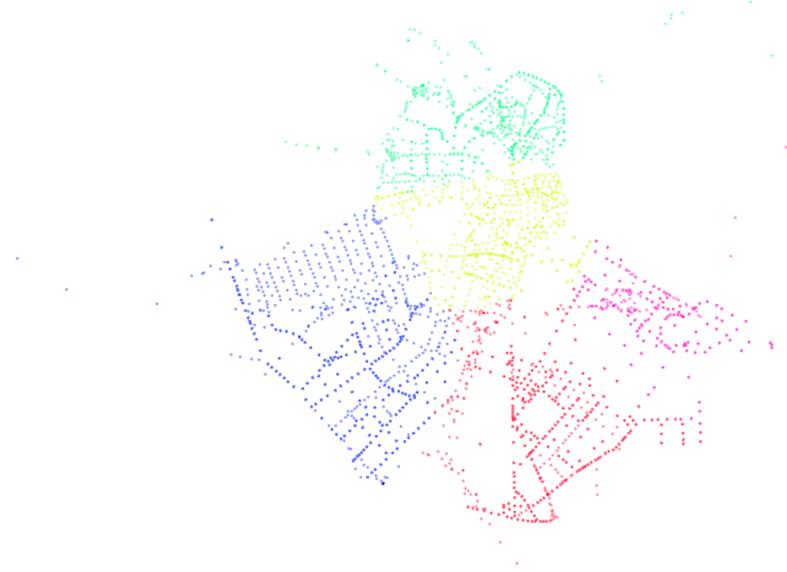
SSCOMP-TOLL



(b) 2 zones derived using SSSCOMP for feature TSI, with $SC=0.194$, $DB=1.252$

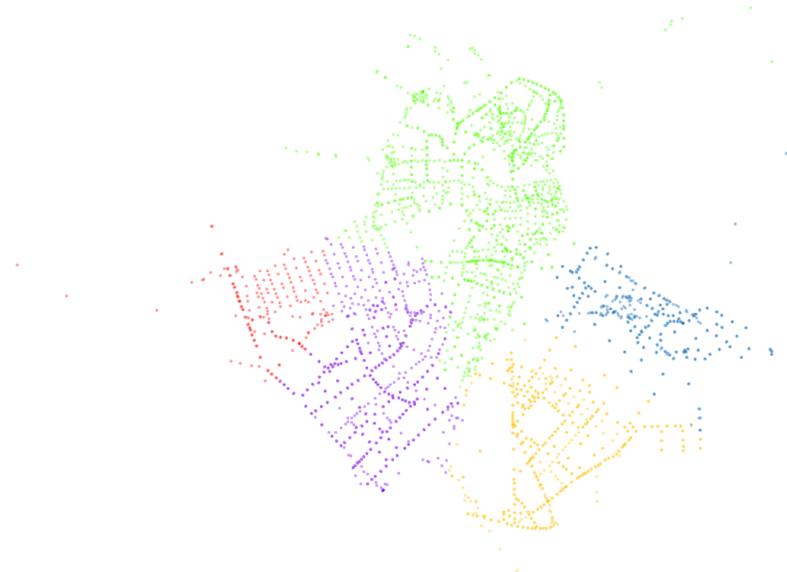
Figure 9: Clustering results and tolling zones (Sparse Subspace Clustering)

OPTICS-TOLL



(a) 5 zones derived using OPTICS for feature TSI, with $SC=0.387$, $DB=0.858$

HDBSCAN-TOLL



(b) 5 zones derived using HDBSCAN* for feature TSI, with $SC=0.400$, $DB=0.776$

Figure 10: Clustering results and Tolling zones (Hierarchical Density-based Clustering)

SSCEL-optTSI



(a) 2 zones derived using SSCEL for feature TSI, with $SC=0.039$, $DB=5.043$

SSCOMP-optTSI



(b) 2 zones derived using SSSCOMP for feature TSI, with $SC=0.102$, $DB=4.009$

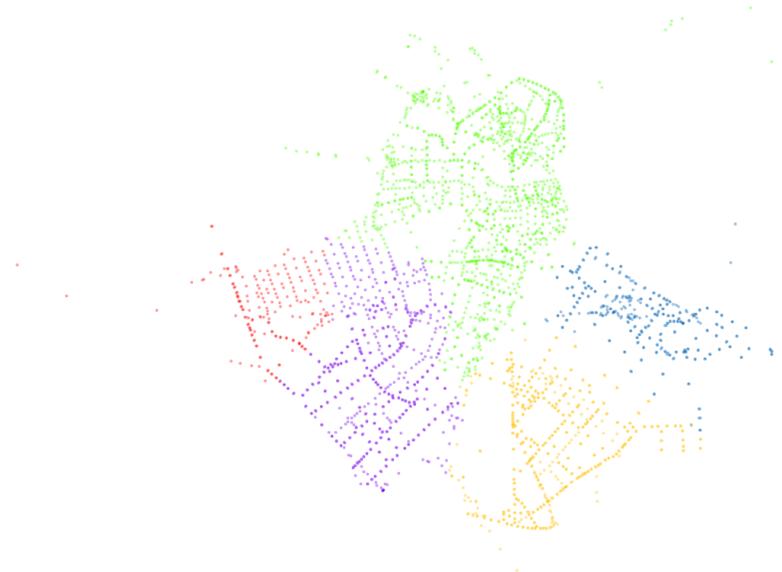
Figure 11: Clustering results and tolling zones (Sparse Subspace Clustering)

OPTICS-optTSI



(a) 5 zones derived using OPTICS for feature TSI, with $SC=0.387$, $DB=0.858$

HDBSCAN-optTSI



(b) 5 zones derived using HDBSCAN* for feature TSI, with $SC=0.400$, $DB=0.776$

Figure 12: Clustering results and Tolling zones (Hierarchical Density-based Clustering)

461 **6. Conclusions and Future Work**

462 In this paper we investigated the use of sparse subspace clustering methods to define tolling zones
463 for distance-based tolling schemes, and their impact on traffic network performance using a predictive
464 real-time distance-based toll optimization framework. Experiments were conducted on the real-world
465 urban network of the Boston Central Business District. We determined that the best network perfor-
466 mance comes from the use of distance-based tolling zones derived from HDBSCAN*, when using Travel
467 Speed Index data. Performance using only 2 tolling zones acquired via the SSCOMP sparse subspace
468 clustering variant was found to be comparable to that of a 5-zone, HDBSCAN*-based derivation, so, in
469 cases where minimizing computational effort is one of the primary objectives, as is the goal of this work,
470 it should be considered as a viable alternative. Despite the fact that all clustering approaches produced
471 tolling zone derivations which, as part of our framework, contributed to significant performance gains,
472 when compared to the No Toll case, we observed large differences in performance between tolling zone
473 derivations acquired via the sparse subspace clustering variants. Specifically, for this particular dataset,
474 the SSCEL variant of sparse subspace clustering produced low quality clustering, due to the low degree
475 of spatial compactness. This warrants further investigation, however, overall, tolling zone derivations
476 acquired from both types of clustering methods, yielded significant benefits on network performance
477 and even outperformed a predictive distance-based tolling scheme that treated the network as a single
478 zone. Finally, the results also underscore the importance of relying not solely on clustering perfor-
479 mance metrics but also the reasonableness of the clusters themselves using context-specific knowledge
480 of demand patterns, land-use etc.

481 In future work, we aim to evaluate alternate clustering methods for systematic tolling zone derivation
482 as part of the distance-based tolling optimization framework. Compared to our No Toll base case,
483 social welfare and network performance results suggest that the clustering can produce distance-based
484 tolling zones with considerable positive impact. We are also in the process of investigating alternative
485 solution approaches, including Bayesian and Surrogate Optimization, and comparing toll optimization
486 framework performance to that of our currently used solution approach.

487 **Author statement**

488 The authors confirm contribution to the paper as follows:

489 **Antonis F. Lentzakis:** Conceptualization, methodology, visualization, investigation, formal analy-
490 sis, writing-Original draft preparation, writing-Review & Editing **Ravi Seshadri:** Conceptualization,
491 methodology, writing-Original draft preparation, writing-Review & Editing **Moshe Ben-Akiva:** Con-
492 ceptualization, methodology, supervision.

493 **Acknowledgements**

494 This research was supported by the National Research Foundation of Singapore through the Singapore-
495 MIT Alliance for Research and Technology's FM IRG research programme.

496 **References**

- 497 Ankerst, M., Breunig, M.M., Kriegel, H.P., Sander, J., 1999. Optics: ordering points to identify the
498 clustering structure, in: ACM Sigmod record, ACM. pp. 49–60.
- 499 Azevedo, C.L., Seshadri, R., Gao, S., Atasoy, B., Akkinpally, A.P., Christofa, E., Zhao, F., Trancik,
500 J., Ben-Akiva, M., 2018. Tripod: sustainable travel incentives with prediction, optimization, and
501 personalization, in: the 97th Annual Meeting of Transportation Research Board.
- 502 Bako, L., 2011. Identification of switched linear systems via sparse optimization. *Automatica* 47,
503 668–677.
- 504 Ben-Akiva, M., Koutsopoulos, H.N., Antoniou, C., Balakrishna, R., 2010. Fundamentals of Traffic
505 Simulation. New York, NY. chapter 10-Traffic Simulation with DynaMIT. International Series in
506 Operations Research and Management Science.
- 507 Bonsall, P.W., Palmer, I.A., 1997. Do time-based road-user charges induce risk-taking? - results from
508 a driving simulator. *Traffic Engineering and Control* 38, 200–203.

509 Campello, R.J., Moulavi, D., Sander, J., 2013. Density-based clustering based on hierarchical density
510 estimates, in: Pacific-Asia conference on knowledge discovery and data mining, Springer. pp. 160–172.

511 Chen, S., 2022. Efficient and Equitable Travel Demand Management Using Price and Quantity Controls.
512 Ph.D. thesis. Massachusetts Institute of Technology.

513 Daganzo, C.F., Lehe, L.J., 2015. Distance-dependent congestion pricing for downtown zones. *Trans-*
514 *portation Research Part B* 75, 91–99.

515 Davies, D.L., Bouldin, D.W., 1979. A cluster separation measure. *IEEE Transactions on Pattern*
516 *Analysis and Machine Intelligence* 1, 224–227.

517 De Palma, A., Kilani, M., Lindsey, R., 2005. Congestion pricing on a road network: A study using the
518 dynamic equilibrium simulator metropolis. *Transportation Research Part A: Policy and Practice* 39,
519 588–611.

520 Elhamifar, E., Vidal, R., 2009. Sparse subspace clustering, in: 2009 IEEE Conference on Computer
521 Vision and Pattern Recognition, IEEE. pp. 2790–2797.

522 Elhamifar, E., Vidal, R., 2013. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE*
523 *transactions on pattern analysis and machine intelligence* 35, 2765–2781.

524 Eliasson, J., Mattsson, L.G., 2006. Equity effects of congestion pricing: quantitative methodology and
525 a case study for stockholm. *Transportation Research Part A: Policy and Practice* 40, 602–620.

526 Geroliminis, N., Levinson, D.M., 2009. Cordon pricing consistent with the physics of overcrowding, in:
527 *Transportation and Traffic Theory 2009: Golden Jubilee*. Springer, pp. 219–240.

528 Gu, Z., Saberi, M., 2019a. A bi-partitioning approach to congestion pattern recognition in a congested
529 monocentric city. *Transportation Research Part C: Emerging Technologies* 109, 305–320.

- 530 Gu, Z., Saberi, M., 2019b. A simulation-based optimization framework for urban congestion pricing
531 considering travelers' departure time rescheduling, in: 2019 IEEE Intelligent Transportation Systems
532 Conference (ITSC), IEEE. pp. 2557–2562.
- 533 Gu, Z., Shafiei, S., Liu, Z., Saberi, M., 2018. Optimal distance-and time-dependent area-based pricing
534 with the network fundamental diagram. *Transportation Research Part C: Emerging Technologies* 95,
535 1–28.
- 536 Gupta, S., Seshadri, R., Atasoy, B., Pereira, F.C., Wang, S., Vu, V.A., Tan, G., Dong, W., Lu, Y.,
537 Antoniou, C., Ben-Akiva, M., 2016. Real time optimization of network control strategies in dynamit2.
538 0, in: *Transportation Research Board 95th Annual Meeting*.
- 539 Gupta, S., Seshadri, R., Atasoy, B., Prakash, A.A., Pereira, F., Tan, G., Ben-Akiva, M., 2020.
540 Real-time predictive control strategy optimization. *Transportation Research Record* doi:[10.1177/
541 0361198120907903](https://doi.org/10.1177/0361198120907903).
- 542 Hashemi, A., Vikalo, H., 2018. Evolutionary self-expressive models for subspace clustering. *IEEE*
543 *Journal of Selected Topics in Signal Processing* 12, 1534–1546.
- 544 Ji, Y., Geroliminis, N., 2012. On the spatial partitioning of urban transportation networks. *Trans-*
545 *portation Research Part B: Methodological* 46, 1639–1656.
- 546 Lee, H., Battle, A., Raina, R., Ng, A.Y., 2007. Efficient sparse coding algorithms, in: *Advances in*
547 *neural information processing systems*, pp. 801–808.
- 548 Lehe, L., 2019. Downtown congestion pricing in practice. *Transportation Research Part C: Emerging*
549 *Technologies* 100, 200–223.
- 550 Lentzakis, A.F., Seshadri, R., Akkinepally, A., Vu, V.A., Ben-Akiva, M., 2020. Hierarchical density-
551 based clustering methods for tolling zone definition and their impact on distance-based toll optimiza-
552 tion. *Transportation Research Part C: Emerging Technologies* 118, 102685.

553 Lentzakis, A.F., Su, R., Wen, C., 2014. Time-dependent partitioning of urban traffic network into
554 homogeneous regions, in: Control Automation Robotics & Vision (ICARCV), 2014 13th International
555 Conference on, IEEE. pp. 535–540.

556 Li, Y., Xiao, J., 2020. Traffic peak period detection using traffic index cloud maps. *Physica A: Statistical*
557 *Mechanics and its Applications* , 124277.

558 Litman, T., 2019. Congestion costing critique-critical evaluation of the “urban mobility report”-9
559 september 2019 .

560 Liu, Z., Wang, S., Meng, Q., 2014. Optimal joint distance and time toll for cordon-based congestion
561 pricing. *Transportation Research Part B* 69, 81–97.

562 LTA, 2016. Tender awarded to develop next generation electronic road pricing system.

563 LTA, 2021. Installation of on-board units for next-generation erp system delayed to 2023 due to global
564 microchip shortage.

565 Lu, L., Xu, Y., Antoniou, C., Ben-Akiva, M., 2015a. An enhanced spsa algorithm for the calibration
566 of dynamic traffic assignment models. *Transportation Research Part C: Emerging Technologies* 51,
567 149–166.

568 Lu, Y., Seshadri, R., Pereira, F., OSullivan, A., Antoniou, C., Ben-Akiva, M., 2015b. Dynamit2.0:
569 Architecture design and preliminary results on real-time data fusion for traffic prediction and crisis
570 management, in: *Proceedings of IEEE 18th International Conference on Intelligent Transportation*
571 *Systems, Spain*. pp. 2250–2255.

572 Meng, Q., Liu, Z., Wang, S., 2012. Optimal distance tolls under congestion pricing and continuously
573 distributed value of time. *Transportation Research Part E: Logistics and Transportation Review* 48,
574 937–957.

- 575 de Palma, A., Lindsey, R., 2011. Traffic congestion pricing methodologies and technologies. *Trans-*
576 *portation Research Part C: Emerging Technologies* 19, 1377–1399.
- 577 Pham, D.S., Budhaditya, S., Phung, D., Venkatesh, S., 2012. Improved subspace clustering via exploita-
578 tion of spatial constraints, in: *2012 IEEE Conference on Computer Vision and Pattern Recognition*,
579 IEEE. pp. 550–557.
- 580 Rao, S., Tron, R., Vidal, R., Ma, Y., 2009. Motion segmentation in the presence of outlying, incomplete,
581 or corrupted trajectories. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 32, 1832–
582 1845.
- 583 Rousseeuw, P.J., 1987. Silhouettes: a graphical aid to the interpretation and validation of cluster
584 analysis. *Journal of computational and applied mathematics* 20, 53–65.
- 585 Saeedmanesh, M., Geroliminis, N., 2017. Dynamic clustering and propagation of congestion in hetero-
586 geneously congested urban traffic networks. *Transportation Research Procedia* 23, 962–979.
- 587 Schubert, E., Sander, J., Ester, M., Kriegel, H.P., Xu, X., 2017. Dbscan revisited, revisited: why and
588 how you should (still) use dbscan. *ACM Transactions on Database Systems (TODS)* 42, 19.
- 589 Simoni, M., Pel, A., Waraich, R., Hoogendoorn, S., 2015. Marginal cost congestion pricing based on the
590 network fundamental diagram. *Transportation Research Part C: Emerging Technologies* 56, 221–238.
- 591 Simoni, M.D., Kockelman, K.M., Gurusurthy, K.M., Bischoff, J., 2019. Congestion pricing in a world of
592 self-driving vehicles: An analysis of different strategies in alternative future scenarios. *Transportation*
593 *Research Part C: Emerging Technologies* 98, 167–185.
- 594 Smith, M.J., May, A.D., Wisten, M.B., Milne, D.S., Van Vliet, D., Ghali, M.O., 1994. A comparison of
595 the network effects of four road-user charging systems. *Traffic Engineering and Control* 35, 311–315.
- 596 Sun, X., Liu, Z., Thompson, R., Bie, Y., Weng, J., Chen, S., 2016. A multi-objective model for cordon-

597 based congestion pricing schemes with nonlinear distance tolls . *Journal of Central South University*
598 23, 1273–1282.

599 Van Den Berg, V., Verhoef, E.T., 2011. Winning or losing from dynamic bottleneck congestion pricing?:
600 The distributional effects of road pricing with heterogeneity in values of time and schedule delay.
601 *Journal of Public Economics* 95, 983–992.

602 Yang, L., Saigal, R., Zhou, H., 2012. Distance-based dynamic pricing strategy for managed toll lanes.
603 *Transportation Research Record: Journal of the Transportation Research Board* 2283, 90–99.

604 You, C., Li, C.G., Robinson, D.P., Vidal, R., 2016a. Oracle based active set algorithm for scalable
605 elastic net subspace clustering, in: *Proceedings of the IEEE conference on computer vision and*
606 *pattern recognition*, pp. 3928–3937.

607 You, C., Robinson, D., Vidal, R., 2016b. Scalable sparse subspace clustering by orthogonal matching
608 pursuit, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
609 3918–3927.

610 Zhang, K., Chen, Y., Nie, Y.M., 2019. Hunting image: Taxi search strategy recognition using sparse
611 subspace clustering. *Transportation Research Part C: Emerging Technologies* 109, 250–266.

612 Zheng, N., R erat, G., Geroliminis, N., 2016. Time-dependent area-based pricing for multimodal systems
613 with heterogeneous users in an agent-based environment. *Transportation Research Part C: Emerging*
614 *Technologies* 62, 133–148.

615 Zheng, N., Waraich, R.A., Axhausen, K.W., Geroliminis, N., 2012. A dynamic cordon pricing scheme
616 combining the macroscopic fundamental diagram and an agent-based traffic model. *Transportation*
617 *Research Part A: Policy and Practice* 46, 1291–1303.

618 Zhu, F., Ukkusuri, S.V., 2015. A reinforcement learning approach for distance based dynamic tolling
619 in the stochastic network environment. *Journal of Advanced Transportation* 49, 247–266.