



# City Research Online

## City, University of London Institutional Repository

---

**Citation:** Guimera Busquets, J., Alonso, E. & Evans, A. (2018). Air itinerary shares estimation using multinomial logit models. *Transportation Planning and Technology*, 41(1), pp. 3-16. doi: 10.1080/03081060.2018.1402742

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

---

**Permanent repository link:** <https://openaccess.city.ac.uk/id/eprint/18606/>

**Link to published version:** <https://doi.org/10.1080/03081060.2018.1402742>

**Copyright:** City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

**Reuse:** Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

---

---



1    **Air Itinerary Shares estimation using Multinomial Logit Models**

2    Ms. Judit Guimera Busquets<sup>1</sup>

3    *Department of Mechanical Engineering & Aeronautics, City, University of London,*  
4    *EC1V 0HB London, United Kingdom*

5    Dr. Eduardo Alonso

6    *Department of Computer Science, City, University of London, EC1V 0HB London,*  
7    *United Kingdom*

8    Dr. Antony D. Evans

9    *Crown Consulting, Inc., Moffett field, CA 94035, USA*

---

<sup>1</sup> Corresponding author: Judit Guimera Busquets

Email: Judit.Guimera-Busquets.2@city.ac.uk

10    **Air Itinerary Shares estimation using Multinomial Logit Models**

11       The main goal of this study is the development of an aggregate air itinerary  
12      market share model. In order to achieve this, multinomial logit models are  
13      applied to distribute the city-pair passenger demand across the available  
14      itineraries. The models are developed at an aggregate level using open-source  
15      booking data for a large group of city-pairs within the US Air Transport System.  
16      Although there is a growing trend in the use of discrete choice models in the  
17      aviation industry, existing air-itinerary share models are mostly focused on  
18      supporting carrier decision-making. Consequently, those studies define itineraries  
19      at a more disaggregate level, using variables describing airlines and time  
20      preferences. In this study, we define itineraries at a more aggregate level, i.e., as a  
21      combination of flight segments between an origin and destination, without further  
22      insight into service preferences. Although results show some potential for this  
23      approach, there are challenges associated with prediction performance and  
24      computational intensity.

25       Keywords: word; air itinerary shares; discrete choice models; multinomial logit;  
26      aggregation level;

27    **1. Introduction**

28      Good forecasts of future demand for air traffic as well as good forecasts of how airlines  
29      are likely to serve this demand are essential to enable supply to adapt to growth in  
30      demand. While the majority of existing research focuses on improving air travel  
31      demand models, there is a growing interest in developing better itinerary share models  
32      than those that already exist. Itinerary share models can be crucial to support airline  
33      network planning and scheduling since important decisions on resources allocation and  
34      pricing are made based on itinerary demand. These decisions are essential as airlines  
35      plan their operations, purchase equipment and make strategic decisions. Airport  
36      authorities also benefit from good forecasts, given the long timescales associated with  
37      airport development and capacity expansion. Improving the accuracy of itinerary share

38 models is therefore a powerful tool for airline and airport authority planning and  
39 decision making, translating into more efficient operations, improved revenue  
40 management and increase profitability. Consequently, for the past 15 years, efforts have  
41 focused on developing this type of model, shifting away from the Quality of Service  
42 indices (QSI) used during the period when the industry was regulated, and/or more  
43 simplistic approaches – such as time-series and simplistic probability models based on  
44 historical trends – (Garrow, 2010). In contrast, discrete choice models model demand by  
45 capturing *how* individuals make decisions and trade-offs among airports, airlines, price,  
46 level of service and other factors that define the air passenger journey.

47 Most of the current research centres on developing innovative approaches using  
48 such discrete choice modelling. These approaches, which aim to model competition and  
49 customer behaviour to determine air-travel itinerary shares (also known as demand  
50 assignment models), are expected to more accurately predict air travel demand. While  
51 most of the discrete choice models applied in urban transport are built using  
52 disaggregate data and include information about the individual making the decision –  
53 i.e. the passenger –; in air transport, data disaggregation as well as data accessibility are  
54 limiting factors. The need to quickly adapt to changes in demand makes flexibility  
55 crucial for carriers and other stakeholders in the industry. For this reason, most of the  
56 models built to support decision-making rely on booking data, which is generally  
57 proprietary. Furthermore, airlines do not typically record much of the passenger data  
58 that is relevant to passenger decision making, such as age, gender and income. This data  
59 is not typically available, except for a small subset of passengers through surveys,  
60 which are time consuming and costly to complete.

61 Most of the early studies on demand assignment for air travel focus on studying  
62 the distribution of demand across one single dimension, i.e. only focusing on modelling

63 passenger choice in terms of one single criteria, such as airport-choice or airline choice.  
64 These early models were mostly applied to analyse air travellers' choice within multi-  
65 airport cities or regions – i.e., airport choice models (Hansen, 1995; Windle & Dreesner,  
66 1995) – or across airlines – airline choice models (Proussaloglou & Koppelman, 1995)  
67 –. Although the former is the most widely studied topic in discrete choice modelling  
68 within air transport, and has given a deeper understanding to the relationship between  
69 airport attributes and airport market share, a more aggregated assignment of air travel  
70 volume is also needed. Only a few studies present approaches for itinerary market share  
71 estimation across multiple dimensions (i.e., modelling a passenger's simultaneous  
72 choice in terms of multiple criteria, e.g., airline, flight time, fare-class etc.) using  
73 discrete choice modelling. Of those, early models used a multinomial logit (MNL)  
74 approach (Adler, 2001; Coldren *et al.*, 2003; Grosche and Rothlauf, 2007; Atasoy and  
75 Bierlaire, 2012), while more recent models also apply nested logit (NL) models  
76 (Coldren and Koppelman, 2005; Hsiao and Hansen, 2011), mixed multinomial logit  
77 (MMNL) models (Warburg *et al.*, 2006) and other alternatives approaches (Gramming  
78 *et al.* 2005; Carrier, 2008). The mentioned aggregate passenger-allocation studies can  
79 be classified according to the type of data they are based on: revealed preference data  
80 (RP) or booking data; stated preferences (SP) data or survey data; or a combination of  
81 both. Studies using RP data do not usually provide full insight into passenger choice  
82 behaviour since models are estimated based on real booking data, and no information  
83 regarding other alternatives at the moment of booking is generally available. This  
84 limitation often leads to RP models performing poorly due to the high demand  
85 inelasticity of the booking data used to estimate the model (Garrow, 2010). In contrast,  
86 SP data collected from surveys allows for modelling of new non-existing alternatives, as  
87 well as more accurate estimation of the sensitivity of travellers to characteristics of their

88 trips. However, studies using SP data may be subject to bias due to the nature of the  
89 experiment in which the individuals are asked to make hypothetical choices by making  
90 trade-offs among the attributes of the choice set (e.g., level of service, air fare etc.).  
91 Although such surveys provide a customer response to a wider range of choices,  
92 providing a better estimate of how individuals make tradeoffs, they are tailored to the  
93 needs of the survey writer, which limits the natural range of choice sets to only those  
94 that the survey writer is aware of (Garrow , 2010; Louviere *et al.*, 1999). Studies based  
95 on SP data are also often limited to a small range of markets, limiting their application  
96 to a small network set.

97       Although the models applied in the studies described above are generally  
98 effective for the purposes to which they are applied, they do not allow for an estimation  
99 of how passenger market demand is distributed across the available itineraries at the  
100 most aggregate level, only considering average market air fare and travel time, level of  
101 service and basic airport attributes as inputs.

102       This paper presents the full air itinerary share model introduced by Busquets *et*  
103 *al.* (2016), refined to better capture passenger choice effects, model validation, and  
104 estimated at the most aggregate level possible, linking annual city-pair demand to the  
105 different itineraries available within the entire US Air Transport System (ATS).

106       The remainder of the paper is structured as follows: The paper's objectives are  
107 presented in Section 2. The modelling approach is detailed in Section 3, with  
108 information regarding the input variables used to estimate the model. The model is  
109 estimated on one dataset, and validated on another. Section 4 provides information  
110 about these two datasets. Modelling results are presented in Section 5, followed by the  
111 model validation results in Section 6 and a discussion on future work in Section 7.

112    **2. Objectives**

113    The primary objective of this research is to develop an air itinerary choice model to  
114    directly estimate the distribution of passenger demand across available routes for a  
115    given O-D pair, using only aggregate data describing average air fare and travel time,  
116    level of service and basic airport characteristics. Ultimately, this model will be  
117    combined with models for forecasting air travel demand and air traffic, all within the  
118    same 3-stage framework (described in Busquets *et al.*, (2015)). This framework consists  
119    of the following stages:

- 120        (1) Forecast city-pair passenger demand;  
121        (2) Distribute this demand across available itineraries; and  
122        (3) Forecast air traffic as a function of route demand.

123            This modelling approach is inspired by previous research that focused on  
124    improving the Federal Aviation Administration's (FAA) forecasting methodology and  
125    for which further potential improvements have been identified. The 3-stage framework  
126    is expected to allow for identification of the key drivers of evolution in the US ATS as  
127    well as to predict future air traffic growth within the US ATS. In order to achieve these  
128    objectives, the approach includes three elements beyond that of the existing research:

- 129        • The use of data mining techniques to model the US ATS evolution in order to  
130            predict air traffic with improved accuracy and precision levels while maintaining  
131            the simplicity if existing econometrics, gravity and time-series models.  
132        • The consideration of a larger set of explanatory variables than is typically  
133            considered in existing air traffic forecasting approaches.  
134        • Explicitly modelling the distribution of city-pair passenger demand between  
135            itineraries.

136 This paper addresses the last of these three elements, which develops the  
137 framework's stage 2 – to distribute passenger demand across available itineraries. The  
138 approach described in this paper is therefore expected to:

- 139 • Highlight the most important factors underlying the air traveller's choice  
140 behaviour within the domestic US ATS;  
141 • Perform air itinerary share model refinement and verification for the entire US  
142 ATS following previously work (Busquets *et al.*, 2016); and  
143 • Explicitly model the distribution of city-pair passenger demand between  
144 itineraries within the US ATS.

145 The model presented in this paper is expected to generate better predictions of airport-  
146 pair air traffic flows once integrated with the air traffic demand model presented by  
147 Busquets *et al.*, (2015).

148 **3. Approach**

149 **Data**

150 Based on the literature review, there are a large number of factors that describe an  
151 itinerary. An itinerary, as defined in this paper, is a flight segment or combination of  
152 flight segments connecting a given city-pair. In this study, itineraries are either non-  
153 stop, or one-stop (i.e., a combination of two flight segments involving an aircraft change  
154 during the connection). Considering constraints in data availability and the different  
155 attributes that are considered to contain the most relevant information for an itinerary,  
156 the input variables for the itinerary market share model are chosen as described in Table  
157 1.

158 [Table 1]

159        The output variable for the model developed in this paper is the market share ( $S_i$ )  
160        of a given itinerary  $i$ . This is defined as the ratio of the demand of the itinerary  $i$  ( $d_i$ ), to  
161        the total demand for the market served by itinerary  $i$  ( $D_m$ ), as shown in Eq. (1). The total  
162        demand for market  $m$  is given by the sum of passengers travelling on all itineraries that  
163        serve that market.

164                  
$$S_i = \frac{d_i}{D_m} \quad (1)$$

165        ***Detailed Forecasting Methodology***

166        Following the work presented by Busquets *et al.* (2015), which introduced the 3-stage  
167        model described in §2 to forecasting future air traffic levels, this paper focuses on fully  
168        developing its stage 2 – to distribute passenger demand across available itineraries. The  
169        objective of this phase is therefore to transform Origin-Destination (O-D) demand by  
170        city-pair into passenger demand by airport-pair using an air itinerary choice model.

171                  Stage 2 of the 3-stage model described by Busquets *et al.* (2015) consists of 2  
172        steps: identification of available itineraries estimated using logistic regression  
173        (described in detail in Busquets *et al.* (2015)), followed by the distribution of the O-D  
174        demand by city-pair obtained from the O-D demand model (stage 1 in the 3-stage model  
175        described by Busquets *et al.* (2015)) across the available itineraries using a discrete  
176        choice model. The first step is motivated by the scope of this research to improve the  
177        current FAA's forecasting methodology while maintaining the simplicity of current  
178        models and is inspired by a previous research (Kotegawa, 2012). The second step is the  
179        focus of this paper. This air itinerary model allows the flight segment passenger demand  
180        by airport-pair to be estimated, based on the passenger itinerary demand from all O-D  
181        city-pairs. It is not feasible to develop a model for each possible O-D market, so in  
182        order to apply the discrete choice model, the US is divided into five regions, as done by

183 Coldren, et al. (2003): four Continental time zones (Central, East, Mountain and West)  
184 and a region for Alaska and Hawaii. This specific O-D market grouping is an attempt to  
185 capture similarities among all city-pairs. The number and nature of these regional  
186 clusters will be modified using clustering techniques in future work. Given these  
187 regions, 18 region-pairs have been defined considering all 16 possible combinations of  
188 the Continental time zones – e.g., Central-Central (C-C), Central-East (C-E), Central-  
189 Mountain (C-M), Central-West (C-W), etc., West-Mountain (W-M), West-West (W-W)  
190 –; as well as a region-pair for Alaska and Hawaii to the Continental US and an region-  
191 pair for the Continental US to Alaska and Hawaii. For each region-pair, henceforth  
192 referred to as an 'entity', an air itinerary share model is developed.

193        This attempts to model the aggregate share of all or groups of decision makers -  
194 i.e., air travellers - choosing each alternative as a function of the trip characteristics. In  
195 contrast to existing research, the itinerary share estimation is done at the most  
196 aggregate level, without considering variables specific to the traveller, such as  
197 passenger preferences and perceptions, or variables specific to the service provider,  
198 such as airline operating the given route, departure time or aircraft type, among others.  
199 Instead, only attributes related to average air fare and travel time, level of service and  
200 basic airport characteristics are considered. The focus of the model is to estimate the  
201 distribution of annual passenger market demand among itineraries, which will be used  
202 as one of the input variables in the third stage of the air traffic estimation model  
203 described in §2, per annum.

204        In order to develop the air itinerary share model, RP data is used, avoiding the  
205 risk of response bias and allowing for the consideration of a much larger network of  
206 city-pairs within the US ATS. The RP data used is 10% ticket survey of booking data  
207 from airlines operating within the US domestic market (BTS-RITA, 2003-2010). The

208 city-pairs considered,  $M$ , are all within the domestic US ATS and are defined by origin  
 209 and destination. The universal choice set,  $C$ , is formed for all possible itineraries within  
 210 the entire ATS connecting these city pairs. The choice problem is defined for each city-  
 211 pair,  $m \in M$ , with the choice set being all the possible itineraries connecting that given  
 212 city-pair, represented by  $I_m$ . Each itinerary  $i$  is characterised by a set of attributes such  
 213 as level of service, price, time and basic airport characteristics. As a simplification, only  
 214 two possible levels of service are considered, non-stop and one-stop flights. For the one-  
 215 stop flights, the connections available are through one of a set of 24 US hub airports  
 216 defined for this study.

217 The annual share of passenger demand assigned to each itinerary between a  
 218 given city-pair is modelled as an aggregate multinomial logit (MNL) function and is  
 219 given by Eq. (2) where  $S_i$  is the passenger share assigned to itinerary  $i$ ,  $V_i$  is the utility  
 220 function or value of itinerary  $i$  and the summation is over all itineraries for a given city-  
 221 pair. The utility function ( $V_i$ ) is a linear function of the explanatory variables and  
 222 assumes that each vector of attributes characterizing an alternative can be reduced to a  
 223 scalar value, which expresses the attractiveness of each alternative. Consequently, it is  
 224 expected that the individual or group of individuals will choose the alternative with the  
 225 highest value, maximizing their utility. Equation (3) shows the general expression for  
 226  $V_i$ , where  $X_i$  is the vector of attributes defining alternative  $i$ ; and  $\beta'$  represents the  
 227 coefficients to be estimated capturing the influence of the corresponding attribute on the  
 228 alternative  $i$  (Atasoy & Bierlaire, 2012).

229 
$$S_i = \frac{\exp(V_i)}{\sum_j \exp(V_j)} \quad (2)$$

230 
$$V_i = \beta' \cdot X_i = \beta_1 \cdot X_{i1} + \beta_2 \cdot X_{i2} + \dots + \beta_k \cdot X_{ik} \quad (3)$$

231 Attributes included in the  $X_i$  vector are described in Table 1 (§3). Some interactions  
232 between the attributes are accounted for by the model. After evaluating several model  
233 specifications, the interactions that define the utilities considered in this paper were  
234 identified as follows:

235 • Accessibility: The interaction between airport accessibility information and  
236 multi-airport city information is accounted for (i.e., the *masORIG* and *masDEST*  
237 variables). Four possible interactions are possible, two regarding the origin  
238 airport and two regarding the destination airport. However, because coefficients  
239 need to be normalised, the coefficients regarding accessibility for origin and  
240 destination airports within cities that are not multi-airport systems are set to 0.

241 • From/to hub variables: The interaction between the hub variables (i.e., whether  
242 the itinerary is from and to a hub, only the origin or destination airport is a hub,  
243 or none of the itinerary airports are hubs) and markets that contain at least one  
244 non-stop itinerary is considered. From/to hub variables are normalised by setting  
245 the variable from and to a hub (i.e., the *hub2hub* variable) to 0.

246 During the estimation of the model, for each city-pair considered, the utility and  
247 likelihood function are computed, with the latter being used to calculate the final  
248 estimated log likelihood.

249 Although all 18 air-itinerary share models have been developed, in this paper  
250 estimated results are only presented for six entities (the entities C-M, M-C, C-W, W-C,  
251 M-W and W-M). Due to issues with computational intensity during the estimation  
252 process for some entities, reduced estimation datasets were generated by sampling a  
253 subset of the total number of city-pairs within the given entity. The size of the reduced  
254 estimation datasets was chosen after evaluating preliminary model estimation results

255 obtained when considering different estimation dataset sizes. Due to the aggregate  
256 nature of the data used in this study and the fact that this data represents only a 10%  
257 sample of real booking data, limiting assumptions are implicitly included when  
258 estimating the model. For example, some itineraries have a very small probability of  
259 occurring, heavily influencing the results obtained for the model estimated as well as its  
260 performance. Moreover, due to the large number of city-pairs considered in the  
261 estimation data and the large number of coefficients to be estimated, the model  
262 estimation becomes computationally too intensive. For these reasons, the data is  
263 reduced to 10 datasets containing information on 100 randomly chosen city-pairs, which  
264 are then each used to estimate the model, reducing the complexity of the problem. The  
265 final estimated model coefficients are computed as the average of the 10 different  
266 models. The performance of each of the entities' air itinerary share model is validated  
267 with data not used for the model estimation. Table 2 reports summary statistics for all  
268 the entities. The set of hub airports varies between entities, as some hubs do not make  
269 sense for some entities for geographical reasons. Table 2 shows the busiest flows in the  
270 US ATS network, i.e., the East Coast corridor (East - East entity), the Central corridor  
271 (Central – Central entity) and between the Central region and East Coast (Central-East  
272 and East-Central entities). A total of 17,200 city-pairs and 104,806 itineraries within the  
273 US ATS network are accounted for in the development of the air itinerary share models.

274 To better understand the results obtained from the air itinerary share model,  
275 indicators such as passengers' 'willingness to pay' can be computed. Value of time  
276 (VOT) is the willingness of passengers to pay for one hour of travel and is defined by  
277 Eq. (4), which is computed for each given itinerary  $i$ . Note that because *Travel Fare*  
278 *Ratio* is a function of the average air fare in the market and *Travel Time Ratio* is a  
279 function of the minimum flight time possible in the market, when computing the utility

280  $V_i$ , average air fare ( $\overline{TF}$ ) and minimum flight time ( $TT_{sh}$ ) are also included in the  
281 formulation of VOT.

282 
$$VOT_i = \frac{\partial V_i / \partial time_i}{\partial V_i / \partial price_i} = \frac{\beta_{FlightTimeRatio}}{\beta_{AirFareRatio}} \cdot \frac{\overline{TF}}{TT_{sh}} \quad (4)$$

283 [Table 2]

284

285 Once the itinerary choice model is estimated using the MNL function, Eq. (1) is  
286 applied to compute the market share of passengers on each itinerary. The estimated  
287 passenger demand per itinerary is then used to compute segment demand – i.e.,  
288 passenger demand per airport-pair – which will ultimately be used as an input for stage  
289 3 of the 3-stage model described in §2, as described in detail by Busquets *et al.* (2015).

290 **4. Application**

291 The models described above are applied to a network of 337 airports within the US  
292 ATS, as used in the Aviation Integrated Modelling (AIM) Project (2006). The choice of  
293 the US air transport network is motivated by improving the current FAA's forecasting  
294 methodology, and by the availability of data. The availability of data for the analysis of  
295 air transport systems can be challenging, with the US being one of the few countries to  
296 provide open source data.

297 The RP data used for this study includes passenger demand data and airfares  
298 extracted from the Airline Origin and Destination Survey (DB1B) (BTS-RITA, 2003-  
299 2010), which contains a 10% sample of airline tickets from reporting carriers. Travel  
300 times and costs are also extracted from BTS-RITA (2003-2010). The air itinerary choice  
301 model is estimated using Biogeme (Bierlaire, 2003). Flight delay information is

302 obtained from the FAA Aviation System Performance Metrics (ASPM) database (FAA,  
303 2007-2010).

304 The RP data considered for estimating the model is from 2007, to be in line with  
305 the period considered when estimating the ultimate 3-stage model described by  
306 Busquets *et al.* (2015). The data used to validate the model is from 2010.

307 Once the model is estimated, it will be applied in future work to estimate the  
308 itinerary shares in the same network of 337 airports into the future. These results will  
309 then be compared to those of the Terminal Area Forecasts (TAF) produced by the FAA.

310 **5. Model Estimation Results**

311 Parameter estimates for the six air itinerary share models mentioned above are reported  
312 in Table 3 below. From the entities shown, parameters for the C-W and W-C entities are  
313 estimated using 10 different folds of 100 randomly selected city-pairs. The estimated  
314 coefficients are averaged to define the final model coefficients. For the C-M, M-C, M-  
315 W and W-M entities, the entire estimation dataset is used to estimate the air itinerary  
316 share model. As Table 2 shows, the C-W and W-C entities have 724 city-pairs and just  
317 over 5,200 itineraries, while the rest of the entities' datasets reported in this paper  
318 contain a much lower number of city-pairs and itineraries, making the estimation  
319 process less computationally intensive.

320 Model performance is described using the likelihood ratio test and rho-squared  
321 parameter ( $\rho^2$ ). The likelihood ratio test provides an evaluation of the entire estimated  
322 model by evaluating whether it is possible to reject the null hypothesis that a more  
323 restricted model (i.e., a model with zero coefficients) is equal to the estimated one. The  
324  $\rho^2$  metric is an indicator of overall goodness of fit.

325 All estimated coefficients are statistically significant at the 95<sup>th</sup> percentile  
326 confidence level.

327 The *Travel Fare Ratio* and *Travel Time Ratio* coefficients are both of the  
328 expected sign, negative, indicating that fares and travel time are a resistance to travel. In  
329 contrast, some of the coefficients associated with delay at the origin and destination  
330 airports are positive, suggesting a correlation between delay and itinerary attractiveness,  
331 which is unexpected. For entities C-M, M-C, M-W and W-M, the sign of the  
332 coefficients alternates between positive and negative, indicating a positive correlation  
333 between delay and itinerary attractiveness associated with Mountain (M) airports. For  
334 the C-W and W-C entities both delay parameters are positive. These results may be an  
335 indication of airport importance since larger and/or hub airports are expected to have  
336 more passengers and flights, and therefore higher delay. This suggests that passengers  
337 are more inclined to travel to and from large airports, which is likely because of the  
338 increased number of routing alternatives available at these airports.

339 The coefficients associated with airport accessibility are also positive, with the  
340 exception of the *AccessDEMas* coefficient for the C-M entity and the *AccessORMas*  
341 coefficient for the W-C entity. This is opposite to what one would expect since an  
342 increased travel time to/from an airport is a resistance to air travel, and given the  
343 influence on door-to-door travel time, a negative sign would be expected. However, the  
344 coefficients associated with all airport accessibility time variables are small, - with the  
345 exception of the *AccessORMas* coefficients for the M-C and M-W entities -, indicating  
346 low influence of passenger preferences on itinerary choice.

347 [Table 3]  
348

349        The estimated *Airline Ratio* coefficients tend to be in the order of 10e-2 and  
350        positive - with the exception of the coefficient associated with the C-W entity -,  
351        indicating low influence of passenger preference on itinerary choice. Coefficients  
352        associated with level of service are represented by dummy variables in the models and  
353        are characteristics of every entity. These variables show the passengers preference in  
354        terms of level of service and connecting hub choice. Due to the fact that each entity has  
355        a specific set of hubs and different assumptions have been made in building the  
356        connection alternatives, a comparison of the estimated coefficients across entities is not  
357        possible.

358        For the variables associated with origin and destination hub information (*Ihub*  
359        and *no\_hub*), both coefficients are generally negative, except for the C-W and W-C  
360        entities. One would expect a negative correlation between itinerary attractiveness and  
361        traveling from or to a hub airport (i.e., *Ihub*=1), and also between itinerary  
362        attractiveness and travelling from and to a non-hub airport (i.e., *no\_hub*=1). In both  
363        cases fewer alternatives would exist than for an itinerary between two hubs. The  
364        positive correlation for entities C-W and W-C may be because these sets of variables  
365        interact only with itineraries belonging to markets in which non-stop options exist, and  
366        itineraries from or/and to a non-hub airport may be associated with lower delay as well  
367        as lower travel fare ratio than itineraries from and to a hub.

368        Regarding the model performance, both the likelihood ratio test and rho-squared  
369        parameters for the six entities show reasonable goodness of fit. Although all the models  
370        show a likelihood ratio test large enough to reject the null hypothesis that all  
371        coefficients are equal to zero; rho-squared values tend to be largest for those models for  
372        which the entire dataset has been used during estimation. While the C-M, M-C, C-W  
373        and W-C entities have a rho-squared value of about 0.7; the rho-squared values for the

374 C-W and W-C entities are lower than 0.6. The same trend is found for the other air  
375 itinerary models estimated.

376 To further analyze the results and understand the effect that the level of service  
377 has on the willingness to pay, VOT is computed – using Eq. (4) – for an example case.  
378 Table 4 shows the VOT values for the six air itinerary share models presented in this  
379 paper. For each of the entities an example case has been chosen and the corresponding  
380 VOT value has been computed. Considering that VOT values in the literature are  
381 typically under \$100/hour (Hsiao & Hansen, 2011; Atasoy & Bierlaire, 2012) several  
382 observations can be highlighted from the results presented in Table 4. While the  
383 estimated VOT for the specified city-pair belonging to the W-M entity is high compared  
384 to the literature (i.e., \$144.42/hr), the estimated values for the case examples from the  
385 other entities are well below \$100/hr, and therefore comparable to those found in the  
386 literature. This may be because of a lack of differentiation between fare classes, the  
387 level of aggregation of the data used or the differences between the entities' estimation  
388 datasets.

389 [Table 4]

## 390 **6. Model Results Validation**

391 The estimated air itinerary share models are validated using data associated with city-  
392 pairs existing in the corresponding entity for the first quarter of 2010. To evaluate the  
393 performance of the model, the market share by itinerary predicted by the model is  
394 compared to the observed market share obtained directly from the DB1B dataset (BTS-  
395 RITA, 2003-2010). Absolute errors are averaged across itineraries, shown in Table 5.  
396 Validation results obtained show an average mean absolute error, expressed in terms of  
397 percentage deviation, of 14.2%, ranging from 7.5% for the W-E entity to 27.2% for the

398 M-M entity. Most of the percentage errors in itinerary share are lower than those in the  
399 literature (e.g., the model developed by Coldren et al. (2003) for 2010 passenger  
400 itinerary shares has a mean absolute error of 16.6%). Only the percentage errors  
401 associated with the M-M entity, the Hawaii & Alaska-US Continental entity and US  
402 Continental-Hawaii & Alaska entity are larger. The model specifications and data  
403 aggregation, however, differ markedly, so such a direct comparison of model  
404 performance is difficult.

405 It is believed that the primary differences lie in the fact of the estimation dataset  
406 used to estimate the M-M air itinerary share model has the smallest number of  
407 observations compared to the other entities, as shown in Table 2. The high mean  
408 absolute error values obtained for the Hawaii & Alaska-US Continental entity and the  
409 US Continental-Hawaii & Alaska entity, may be due to the different assumptions  
410 implicit in the datasets. While the rest of the entities contain city-pairs with the same  
411 time-zone difference, these two entities contain a variety of time zones, which may  
412 affect the estimation results.

413 [Table 5.]

## 414 **7. Conclusion and Future Work**

415 In this paper a step is made to improve on existing air traffic forecasting methodologies  
416 through a better understanding of the factors driving demand, supply and network  
417 dynamics. In order to achieve this, an aggregate air itinerary share model is presented  
418 that only uses aggregate data, without further insight into service preferences, in  
419 contrast to other models in the literature. Given this aggregate input data, the developed  
420 model attempts to model demand effects and passenger travel decision more accurately  
421 than is possible using other methods. Ultimately, when integrated into a 3-stage model

422 for air traffic forecasting, better predictions of airport-pair traffic flows are expected.

423 An aggregate multinomial logit model is estimated to predict how market  
424 demand is distributed across available itineraries. In an attempt to capture similarities  
425 between city-pairs, eighteen models are developed, each modelling traffic-flow between  
426 two major regions of the US ATS. In this paper, results for six entities are presented (C-  
427 M, M-C, C-W, W-C, M-W and W-M entities). Due to computational limitations some  
428 of the models are estimated using a reduced dataset containing information about 100  
429 city-pairs in each of 10 runs. Results obtained from the estimated model show high  
430 goodness of fit. All estimated coefficients are significant at the 95th percentile  
431 confidence level and are generally of the expected sign.

432 The estimated models are validated by computing the mean absolute error  
433 between the predicted market share and the observed market share. Data for city-pairs  
434 from the 1st quarter of 2010 is used for validation. Validation results show an average  
435 mean absolute error of 14.2%, ranging from 7.5% for the W-E entity to 27.2% for the  
436 M-M entity. In general, the validation results obtained are slightly better than  
437 comparable numbers in the literature (Coldren et al., 2003). However, because of  
438 differences in model specifications and data aggregation, a direct comparison is  
439 difficult. Model evaluation parameters including likelihood ratio test and Rho-squared  
440 show reasonable values, with the likelihood ratio test values large enough to reject the  
441 null hypothesis and the Rho-squared values showing a reasonable goodness of fit.  
442 Estimated VOTs are found to be in line with those in the literature for all the entities, -  
443 i.e. under \$100/hr -, with the exception of VOT for the W-M entity. This may be  
444 because of a lack of differentiation between fare classes, the level of aggregation of the  
445 data used or the differences between the entities' estimation datasets.

446 Model estimation results obtained to date look promising, showing that the  
447 application of multinomial logit modelling for air itinerary share estimation at the  
448 aggregate level is possible. However, computational intensity is a significant problem,  
449 requiring the approach to be adjusted to estimate the model with reduced datasets of 100  
450 city-pairs in each of 10 runs. This leads to some issues with the estimated coefficients,  
451 and may reduce model performance. Hence, further work will focus on improving  
452 model estimation results through the use of alternative techniques. Those under  
453 consideration include neural networks using various learning algorithms such as  
454 backpropagation and Levenberg-Marquardt.

455 In future work the best performing model will be used to estimate the air  
456 itinerary shares between city-pairs, so that passenger demand by airport-pair can be  
457 predicted and ultimately used as one of the input variables for the final stage of the 3-  
458 stage model. Additionally, by providing more accurate itinerary shares, this model  
459 could be used to aid the decision making process across multiple stakeholders (e.g.  
460 airlines, airport providers, government' agencies, etc.). Route network expansion,  
461 equipment purchase or airport expansion are some examples in which its application  
462 could be beneficial. Moreover, subject to adequate model refinement, there is the  
463 potential of a broader model application to include other transport modes as one of the  
464 choice criteria. This would allow for the analysis of, e.g., competition between air and  
465 ground transport over short distances.

## 466 **Acknowledgements**

467 The authors would like to gratefully acknowledge Dr. Lynnette M. Dray from University  
468 College London, Dr. Bilge Atasoy from Massachusetts Institute of Technology and Dr. Gregory  
469 Coldren from Coldren Choice Consulting Ltd. for their advice on data sources and approach.

471 Table 1. Input variables considered to influence air itinerary market share.

Variable	Name	Description
Level of service	$LoS$	Dummy variable indicating the level of service of the itinerary $i$ (non-stop or one-stop) with respect the best level of service within its market (either non-stop or one-stop with the best connection).
Travel Time Ratio	$TT_i^{Ratio}$	Ratio between travel time of itinerary $i$ and travel time of shortest itinerary in the market $sh$ .
Travel Fare Ratio	$TF_i^{Ratio}$	Average fare paid on a specific itinerary $i$ divided by the market average fare.
Multi-airport system (MAS)	$masORIG_i$	Dummy variable indicating whether the Origin airport is within a multi-airport system or not.
Origin		
Multi-airport system (MAS)	$masDEST_i$	Dummy variable indicating whether the Destination airport is within a multi-airport system or not.
Destination		
Origin airport	$\bar{Dly}_{ORIG}$	Average departure delay of origin airport for the average delay
Destination airport	$\bar{Dly}_{DEST}$	Average arrival delay of destination airport for the average delay
Origin airport	$Access_{ORIG}$	Average distance between city center and origin

---

Accessibility		airport.
Destination airport	$Access_{DEST}$	Average distance between city center and destination
Accessibility		airport.
Origin and destination airports are hubs	$hub2hub_i$	Dummy variable indicating whether itinerary $i$ is between two hub airports.
Either the origin or destination airport is a hub	$1hub_i$	Dummy variable indicating whether itinerary $i$ is from or to a hub airport.
Neither origin nor destination airports are a hub	$no\_hub_i$	Dummy variable indicating whether itinerary $i$ is not from nor to a hub airport.
Airlines Ratio	$AirlinesRatio_i$	Ratio between the number of airlines serving itinerary $i$ and the number of airlines serving the shortest itinerary $sh$ .

---

472

473 Table 2. Summary statistics for all entities.

Origin Region	Destination Region	City-pairs	Itineraries	Nº itineraries per city-pair	Nº Hubs
Hawaii & Alaska	US Continental	438	2,063	19	15
US Continental	Hawaii & Alaska	437	2,052	19	15
Central	Central	1,547	6,335	16	11
Central	East	2,562	14,415	27	19

Central	Mountain	462	1,867	17	17
Central	West	724	5,216	38	19
East	Central	2,552	15,150	38	18
East	East	3,520	21,157	27	17
East	Mountain	508	2,895	18	20
East	West	867	9,268	87	24
Mountain	Central	463	1,899	15	18
Mountain	East	527	3,150	24	18
Mountain	Mountain	134	359	5	6
Mountain	West	252	1,230	27	11
West	Central	724	5,222	38	19
West	East	862	9,274	90	24
West	Mountain	265	1,313	29	11
West	West	356	1,941	31	9
<b>Total</b>		17,200	104,806		

474

475 Table 3. Estimated coefficients for the air itinerary choice model corresponding to  
 476 entities C-M, M-C, C-W, W-C, M-W and W-M.

Variable Name	C - M	M - C	C - W	W - C	M - W	W - M
<i>Level of Service (relevant to every entity)</i>						
Markets Containing Non-Stop itineraries:	--	--	--	--	--	--

	<i>hub2hub</i>	0.000	0.000	0.000	0.000	0.000	0.000
	<i>1hub</i>	-1.590	-2.090	0.013	0.626	-1.090	-0.846
	<i>no_hub</i>	-1.410	-2.650	0.095	0.928	-1.970	-1.380
	<i>Airlines Ratio</i>	0.012	0.017	-0.550	0.017	0.010	0.023
	<i>Travel Fare Ratio (TF<sup>Ratio</sup>)</i>	-3.840	-4.080	-0.789	-1.321	-1.970	-0.754
	<i>Travel Time Ratio (TT<sup>Ratio</sup>)</i>	-1.030	-1.020	-0.170	-0.329	-0.844	-1.530
	$\overline{Dly}_{ORIG}$	-0.174	3.950	0.086	0.627	2.010	-0.026
	$\overline{Dly}_{DEST}$	2.930	-0.092	0.542	0.227	-0.067	1.110
	<i>AccessDEmas</i>	-0.919	0.020	0.044	0.015	0.002	0.331
	<i>AccessORmas</i>	0.098	0.864	0.098	-0.001	0.749	0.005
	<i>LogLikelihood Ratio Test</i>	523,121	435,323	1,030,223	191,252	908,296	906,390
	<i>Rho-squared (<math>\rho^2</math>)</i>	0.724	0.691	0.587	0.559	0.715	0.714

\*Note: All variables are statistically significant at the 95% confidence level.

477

478 Table 4. Comparison between Value of Time for the C-M, M-C, C-W, W-C, M-W and  
479 W-M entities.

Entity	Origin City	Destination City	$\overline{TF}$ (\$)	$TT_{sh}$ (hr)	$VOT$ (\$/hr)
C – M	Chicago	Denver	137.1	2.51	14.66
M – C	Denver	Chicago	136.6	2.24	15.26
C – W	Chicago	Reno	183.5	4.04	9.79
W – C	Reno	Chicago	184.3	3.59	12.78
M – W	Denver	Los Angeles	150.5	2.17	29.76
W – M	Los Angeles	Denver	151.0	2.12	144.42

480

481 Table 5. Mean absolute error in itinerary share computed in terms of percentage  
 482 deviation.

<b>Origin Region</b>	<b>Destination Region</b>	<b>Number of City-pairs</b>	<b>Number of Itineraries</b>	<b>Mean absolute Error in Itinerary Share (%)</b>
Hawaii & Alaska	US Continental	422	1,889	22.60
US Continental	Hawaii & Alaska	435	1,963	24.17
Central	Central	1,490	6,088	13.46
Central	East	2,460	13,457	11.35
Central	Mountain	463	1,931	17.94
Central	West	679	4,814	9.03
East	Central	2,461	13,748	11.14
East	East	3,503	19,487	11.07
East	Mountain	523	3,066	14.06
East	West	785	7,622	8.63
Mountain	Central	464	1,895	16.69
Mountain	East	517	3,049	14.53
Mountain	Mountain	121	309	27.21
Mountain	West	250	1,130	13.22
West	Central	683	4,868	9.40
West	East	786	7,577	7.49
West	Mountain	262	1,243	11.42
West	West	343	1,653	11.97

483

484

485 **References**

486

- 487 Adler, N. (2001). Competition in a deregulated air transportation market. *European  
 488 Journal of Operational Research*, 129(2), 337-345.

- 489 AIM. (2006). *Aviation Integrated Modelling Project*. Retrieved August 15, 2016, from  
490 Aviation Integrated Modelling Project: <http://www.aimproject.aero/>.
- 491 Atasoy, B., & Bierlaire, M. (2012). *An itinerary choice model based on a mixed RP/SP*  
492 *dataset*. ENAC, EPFL. Lausanne: Transport and Mobility Laboratory, ENAC.
- 493 Bierlaire, M. (2003). Biogeme: A free package for the estimation of discrete choice  
494 models. *3rd Swiss Transportation research conference*. Ascona, Switzerland.  
495 Retrieved from Biogeme: A free package for the estimation of discrete choice  
496 models.
- 497 BTS-RITA. (2003-2010). *Bureau of Transportation Statistics - Research and*  
498 *Innovative Technology Administration*. Retrieved September 14, 2014, from  
499 Origin and Destination Survey: DB1B Market for 1003 to 2007:  
500 [http://www.transtats.bts.gov/DL\\_SelectFields.asp?Table\\_ID=247&DB\\_Short\\_N](http://www.transtats.bts.gov/DL_SelectFields.asp?Table_ID=247&DB_Short_N)  
501 ame=Origin%20and%20Destination%20Survey
- 502 Busquets, J. G., Alonso, E., & Evans, A. D. (2015). Application of Data Mining in Air  
503 Traffic Forecasting. *AIAA Aviation Technology, integration and Operations*  
504 *Conference*. Dallas: AIAA.
- 505 Busquets, J. G., Alonso, E., & Evans, A. D. (2016). Predicting Aggregate Air Itinerary  
506 Shares Using Discrete Choice Modeling. *AIAA, Aviation Technology,*  
507 *Integration and Operations Conference*. Washington D.C.
- 508 Carrier, E. (2008). *Modeling the Choice of an Airline Itinerary and Fare Product using*  
509 *booking and Seat Availability Data*. Cambridge, Massachusetts: Massachusetts  
510 Institute of Technology.
- 511 Coldren, G. M., & Koppelman, F. S. (2005). Modeling the competition among air-travel  
512 itinerary shares. GEV model development. *Transportation Research Part A: Policy and Practice*, 39(4), 345-365.

- 514 Coldren, G. M., Koppelman, F. S., Kasturirangan, K., & Mukherjee, A. (2003).
- 515 Modeling aggregate air travel itinerary shares: logit model development at a  
516 major US airline. *Journal of Air Transport Management*, 9(6), 361-369.
- 517 FAA. (2007-2010). *Federal Aviation Administration*. Retrieved September 20, 2015,
- 518 from FAA Operations & Performance Data: <https://aspm.faa.gov>
- 519 FAA. (2014). *Federal Aviation Administration*. Retrieved October 10, 2015, from  
520 Aviation System Performance Metrics (ASPM) Manuals: <https://aspm.faa.gov/>
- 521 Garrow, L. A. (2010). *Discrete Choice Modelling And Air Travel Demand: theory*.
- 522 Aldershot, United Kingdom: Ahsgate.
- 523 Gramming, J., Hujer, R., & Scheidler, M. (2005). Discrete choice modelling in airline  
524 network management. *Journal of Applied Economics*, 20, 467-486.
- 525 Grosche, T., & Rothlauf, F. (2007). *Air Travel Itinerary Market Share Estimation*.
- 526 Manheim: University of Manheim.
- 527 Hansen, M. (1995). Positive feedback model of multiple-airport system. *ASCE Journal  
528 of Transportation Engineering*, 121 (6), 453-460.
- 529 Hsiao, C., & Hansen, M. (2011). A passanger demand model for air transportation in a  
530 hub-and-spoke network. *Transportation Research Part E: Logistics and  
531 Transportation Review*(47), 1112-1125.
- 532 Kotegawa, T. (2012). *Analyzing the evolutionary mechanism of the air transportation  
533 system-of-system using network theory and machine learning algorithm*. West  
534 Lafayette, Indiana: Faculty of Purdue University.
- 535 Louviere, J. J., Meyer, R. J., Bunch, D. S., Carson, R., & Dellaert, B. (1999). Combinig  
536 sources of preference data for modeling complex decision provesses. *Marketting  
537 Letters* 10, 205-2017.

- 538 Proussaloglou, K., & Koppelman, F. S. (1995). Air carrier demand: an analysis of  
539 market share determinants. *Transportation*, 22(4), 371-388.
- 540 Warburg, V., Bhat, C., & Adler, T. (2006). Modeling demographic and unobserved  
541 heterogeneity in air passengers' sensitivity to service attributes in itinerary  
542 choice. *Transportation Research Record: Journal of the Transportation  
Research Board*, 1951, 7-16.
- 543 Windle, R., & Dreesner, M. (1995). Airport Choice in multiple-airport regions. *ASCE  
Journal of Transportation Engineering*, 121(4), 332-337.
- 546
- 547

