



Technologies for Safe & Efficient Transportation

THE NATIONAL USDOT UNIVERSITY
TRANSPORTATION CENTER FOR SAFETY

Carnegie Mellon University

UNIVERSITY of PENNSYLVANIA

Modeling Transit Patterns Via Mobile App Logs

FINAL RESEARCH REPORT

Anthony Tomasic (PI), Aaron Steinfeld (CO-PI), John
Zimmerman (CO-PI), Afsaneh Doryab (CO-PI)

Contract No. DTRT12GUTG11

DISCLAIMER

The contents of this report reflect the views of the authors, who are responsible for the facts and the accuracy of the information presented herein. This document is disseminated under the sponsorship of the U.S. Department of Transportation's University Transportation Centers Program, in the interest of information exchange. The U.S. Government assumes no liability for the contents or use thereof.

The Problem

Transit planners need detailed information of the trips people take using public transit in order to design more optimal routes, address new construction projects, and address the constantly changing needs of a city and metro region. Better transit plans lead to better service and lower costs. Unfortunately, good rider origin-destination information is almost universally unavailable.

In this project we have developed a new method for inferring rider origin-destination (O-D) trip stops in support of transit planning. The meteoric adoption of smartphones along with the growth of transit apps that provide vehicle arrival information at a stop generates a new data resource. Every time a user requests arrival information, the mobile service logs the user's location, the time, and the specific stop they requested information about. Over time, a user's request history functions as "bread crumbs" revealing where and when they have travelled.

The goal of this project is to develop machine-learning models that can infer O-D for a transit service based on the request logs of individual users of mobile transit apps. This project builds on already deployed and extensively used Tiramisu app. In addition to the request log, Tiramisu data includes O-D trips recorded by users that we can use as ground truth for training the machine learning models. We will use this data to build a transit model that can derive results based on model phone app usage. Thus, we can produce models of transit use at a fraction of the cost. This approach also allows continuous O-D modeling, unlike traditional survey and sampling techniques. Note that, as far as we know, the Tiramisu app is a unique source of exact, large-scale, O-D information collected for research purposes. Other researchers have collected O-D using smartphones in small studies, but not through an extensively deployed app with over four years of historical data.

Approach

Current methods for understanding origin-destination and understanding the pinch points in a current transit plan include (1) surveys of riders, (2) interviews with drivers, (3) rider requests and complaints, (4) fare-box data, and (5) automated passenger counter (APC) data. These sources can reveal where additional capacity might help reduce overload, but they cannot predict a more optimal layout of routes and scheduling. While some rail systems attempt to derive origin-destination based on fare-card data since riders check in and out for each trip, the challenge for redesign of train routing is quite different as it mostly requires the construction of new tracks, a prospect that is considerably slower and more expensive than rerouting buses across existing roadways.

The project attempts to generate a generic model that takes Tiramisu mobile app data as input and outputs a highly accurate and generic travel model for the transit community using the app. The approach adopted in this project combines raw tiramisu data with common sense assumptions to address questions about commuter behavior. We build statistical models and provide visualization of commuter behavior, which helps identify

common behavioral patterns, inefficient routes, under served routes and predict the likely destination.

Methodology

We use the log of a user's interaction with the Tiramisu app as a basis of a series of models. Using this data we construct four models:

1. Based on ground-truth O-D information, we build a model that predicts generic travel for users of a transit system (for example, neighborhood to neighborhood). We test this model by comparing it to withheld test data.
2. Based on user lookups, we build a model that translates user lookups into O-D information. For example, suppose a user at 9 am, located in Squirrel Hill, looks up the 61 C (which goes to CMU at 9:10 am). In the afternoon, at 3 pm, the same user, located at CMU, looks up the 61 C going in the reverse direction. In this idealized case, it is reasonable to model this interaction as two trips: Squirrel Hill (9 am) - CMU (9:30 am) and CMU (3:10 pm) to Squirrel Hill (3:30 pm). Of course, this ideal situation rarely occurs, since users look up many different stops and take buses to many locations. We will test this model by comparing predicted O-D to actual O-D.
3. We combine Models 1 and 2 by pipelining the results of Model 2 into Model 1 to predict generic travel based on lookups. We test this model by comparing predicted O-D to actual O-D in withheld test data. In effect, we are measuring the "accuracy lost" by pipelining from the first model to the second model. In addition, we are measuring the overall accuracy of the model.
4. If Model 3 is highly accurate, it is then possible to explore other transit issues via simulation. However, sizeable error becomes problematic since simulating results at scale increases errors in the predictions.

From our analysis, we have derived a simple Model 1. We are currently working on improving this model. Our chief hypotheses for Model 2 are (i) users actually look up arrival times at stops for the trips that they actually take and (ii) users exhibit O-D patterns in look ups. We will test these hypotheses as part of our model construction and evaluate potential for Models 3 and 4.

Findings

Visualization of Origin and Destinations across Pittsburgh

Using mobile app logs we were able to visual various commuter behaviors. We found this representation especially useful in observing and confirming different behavior patterns. A simple heat map can show which stops are the most common origins and destinations in Pittsburgh.

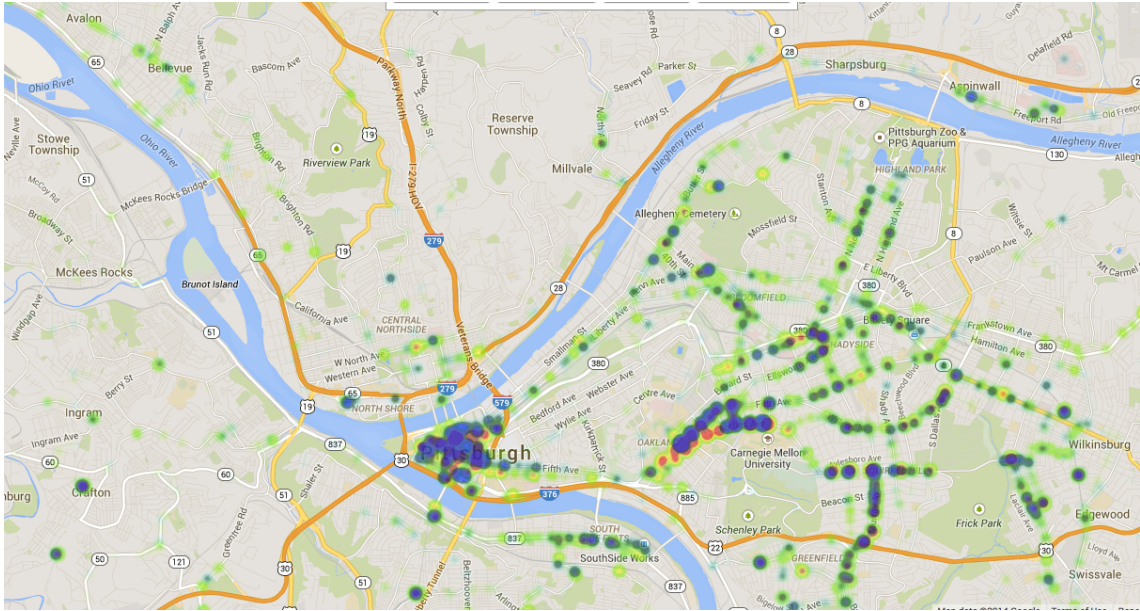


Figure 1. O-D heat-map depicting the stops logged by Tiramisu users in Pittsburgh, red in this figure represents the stop/area that is more of an origin than destination, blue implies it is more of a destination than origin while purple is in between.

Heat maps can help in further analysis. For example, the average weekday travel pattern can be compared with weekend travel pattern in the city.

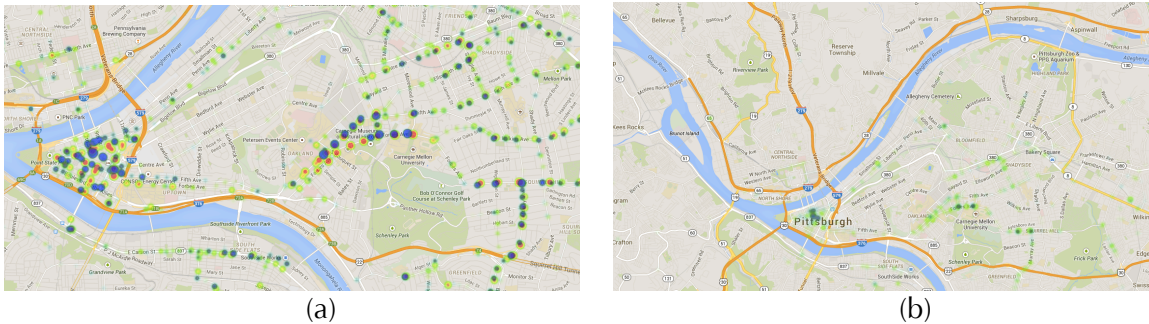


Figure 2. O-D heat-map comparing weekday stop access (a) with weekend stop access (b)

Other visualizations can help identify daily usage patterns, for example by observing the logs of the Tiramisu app we can identify ridership pattern over the course of a day.

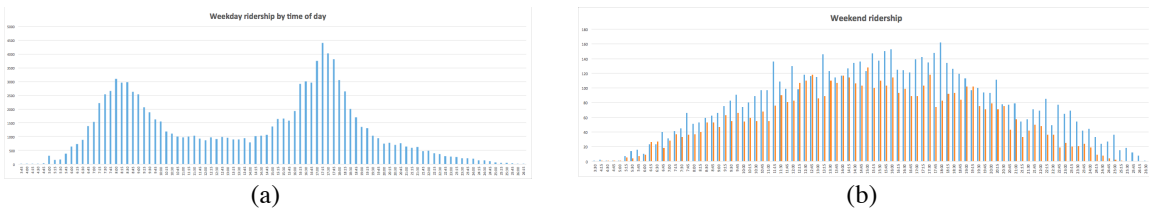


Figure 3. Average ridership as a function of the time of day: (a) Weekday ridership (b) Weekend ridership

Identifying Inefficient Trips

Tiramisu app logs can help identify routes that are frequently used by riders but are not served directly by the transit agency, forcing the commuters to make inefficient transfers during a trip. In the logs if two trips by the same user have timestamps that fall within an hour of one another and destination of the first trip is origin of the second trip but the destination of second trip is not the initial starting point then such trips are considered inefficient. If such trips are being taken by a large number of riders then it may identify a new route that should be served directly. Figure 4, below identifies such routes in Pittsburgh.

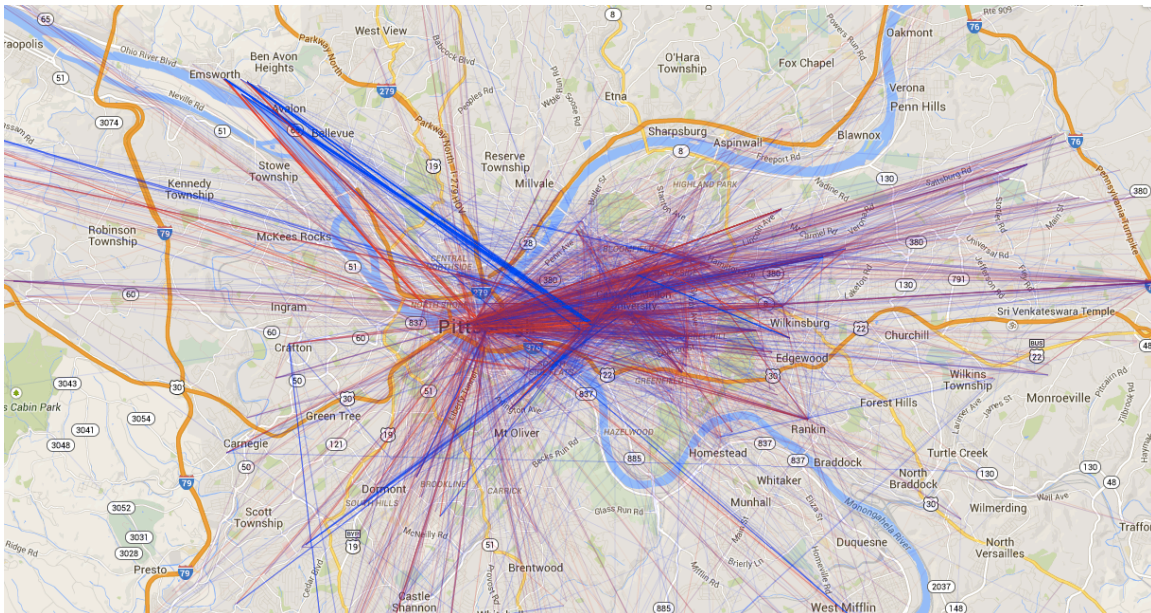


Figure 4. Red lines connect the real origin and destination while the blue lines join the starting point with the final destination and signify the efficient route.

Predicting Commuter Trips

Tiramisu app logs have also been effective in predicting commuter behavior. Given a user's app use history, we can try to predict a trip that is in a user's recorded trips. Simply using the most commonly used stops for origin and destination results in a poor model that is only accurate 3.8% of the time, by considering nearby stops in a square shaped territory 0.5 x 0.5km improves the accuracy to into the 40% range. Raising the dimensions to 1 x 1km and manipulating various input data features can achieve results in the 70-80% accuracy range. This is clearly not high enough to support prediction in large-scale simulations since the error term would grow dramatically. Therefore, it seems our approach is missing key data needed for higher accuracy performance.

Boarding/Exiting Vehicles

In an ideal world, performance on the above O-D estimations and predictions would have higher performance. To address this gap, we considered what new data would lead to better modeling and prediction. Initial explorations suggested that knowledge of when a

person is entering or exiting a vehicle would dramatically improve estimation of ground truth, and therefore overall performance. Therefore, the team began an effort to automatically capture boarding and exiting via smartphone sensors.

Initial explorations of useful sensors included accelerometers, the magnetometer, audio input, and the light sensor. Audio input can be problematic due to the potential for user concerns and possible issues with various wiretap laws. The light sensor may be inappropriate since a user may never take their phone out of their pocket or bag. Past experience with smartphone magnetometers reveals a very noisy signal, which could possibly introduce additional error into a machine learning approach. Therefore, we have focused most of our attention on the accelerometers. Other research teams have had good luck with these sensors.

Our data stream includes contextual features other teams do not have access to. For example, we can combine accelerometer data with the current Tiramisu data stream. We think approach will ultimately lead to better data features for machine learning approaches to O-D modeling, thus leading to better overall predictive performance.

Collaboration

The team has been working with Sean Qian's team to explore how fullness ratings made in Tiramisu align with data collected by Port Authority's automated passenger counter (APC) system. We have identified and extracted a set of data that matches APC data provided to Sean's team by Port Authority. Work on these analyses is underway.

Other Activities

Some of the initial O-D modeling work was done in cooperation with a Masters Capstone project titled, "Activity Recognition Using Mobile Device Sensors." The team included Napat Luevisadpaibul and Wenqing Yuan and was advised by Anthony Tomasic.

The team was active in a variety of stakeholder activities, including regular professional service and meetings with interested industry visitors. For example, Aaron Steinfeld continues to serve on the National Academies of Science, Transportation Research Board, Standing Committee on Accessible Transportation and Mobility (ABE60)¹. He is the Co-Chair of the Technology subcommittee.

Conclusions

Our main finding from this effort is that app logs can provide valuable insight into commuter behavior but additional data streams are needed beyond the data currently captured by Tiramisu if localization under a 1 x 1 km square territory is desired. A 1km square size is useful for some applications, but not sufficient for high precision planning activities.

¹ <https://www.mytrb.org/CommitteeDetails.aspx?CMTID=1164>

Machine learning can also improve with more training data. While we have a lot of data for Pittsburgh, it is spread over a large temporal window that may be introducing some of the error. We will soon advertise Tiramisu in New York City and will re-evaluate some of the models tested in this project if we are able to capture a larger data set.

We foresee that by improving the prediction of rider's destination the application can improve the user interaction and also provide the user relevant information to his/her journey (such as potential detours or stop closures). Eventually, we hope to provide real time data to transit planners for designing better routes and services.

Recommendations drawn from the project:

1. A more formal study of the inefficient stops should be made a part of regular review by the agency's routes review committee. This can be especially beneficial for disabled people. To this end, a new T-SET project starting in 2016 hopes to automate the discovery process of bad connections.
2. We think using sensors on the phone can help detect various activities and events of commuter interest, thus improving the quality of data recorded by the application. We will continue to explore this under other funds.