# Location Segmentation, Inference and Prediction for Anticipatory Computing

Nathan Eagle

MIT Media Laboratory The Santa Fe Institute *nathan@mit.edu*  Aaron Clauset The Santa Fe Institute *aaronc@santafe.edu*  John A. Quinn Makerere University jquinn@cit.mak.ac.ug

#### Abstract

This paper presents an analysis of continuous cellular tower data representing five months of movement from 215 randomly sampled subjects in a major urban city. We demonstrate the potential of existing community detection methodologies to identify salient locations based on the network generated by tower transitions. The tower groupings from these unsupervised clustering techniques are subsequently validated using data from Bluetooth beacons placed in the homes of the subjects. We then use these inferred locations as states within several dynamic Bayesian networks to predict each subject's subsequent movements with over 90% accuracy. We also introduce the X-Factor model, a DBN with a latent variable corresponding to abnormal behavior. We conclude with a description of extensions for this model, such as incorporating additional contextual and temporal variables already being logged by the phones.

#### Introduction

While there has been significant gains in the market of GPSenabled phones, GPS is available on less than 5% of the approximately 4 billion mobile phones in use today. However, every cellular phone has access to information about the proximate cellular towers to which it is attached. The contributions of this paper are to introduce the usage of community structure algorithms to identify salient locations from cellular tower networks, and to build a dynamic, generative model that is extendable to multi-modal data and can be used for behavior prediction. We show that using temporal data from cellular towers, information every phone has access to, it is possible to infer a user's salient locations and anticipate subsequent movements.

We begin by presenting related work that has used continuous cellular tower data for location estimation and location data to train probabilistic graphical models on human movement patterns. We then introduce our own data set, which we believe is the first of its kind to use exclusively randomly sampled subjects. We introduce several segmentation algorithms taken from the community structure literature and apply them to networks of cellular towers. Coupling bluetooth beacon data placed in the homes of each subject with the tower data, we validate the output of the community structure algorithms with the community of towers copresent with the beacon exposures in each subject's home. We then describe several DBNs that use the inferred locations clusters as states to parametrize and predict subsequent movements. One such DBN we use for behavioral modeling includes a latent variable, the X-Factor, corresponding to a binary switch indicative of "normal" or "abnormal" behavior. We conclude with ideas for extensions to these models as future work.

There has recently been a significant amount of research quantifying and modeling human behavior using data from mobile phones. We will highlight a selection of the literature on GSM trace analysis and subsequently discuss recent work on location segmentation and movement prediction from GPS data.

**Cellular Tower Data Analysis** Mobile phones are continuously, passively monitoring signals from proximate cellular towers. However, due to power constraints, a mobile phone typically does not continuously send back similar signals alerting the nearby towers of its particular location. While there has been recent work on analysis of data from mobile phone operators (González, Hidalgo, and Barabási 2008; Onnela et al. 2007), call data records (CDR) from operators only provide estimates of locations when the phone is in use.<sup>1</sup> Therefore, the only method of obtaining continuous cellular tower data is by installing a logging application on the mobile phone itself.

There have been a variety of projects that have involved installing a mobile phone application that logs visible cellular towers and Bluetooth devices on a set of subjects phones including HIIT's Context project, MIT's Reality Mining project (Eagle and Pentland 2006) and the PlaceLab (Chen et al. 2006; LaMarca et al. 2005) research at Intel Research. Additionally, other research projects have demonstrated the utility of cellular tower data for a broad spectrum of applications ranging from contextual image tagging (Davis et al. 2004) to inferring the mobility of an individual (Sohn et al.

Copyright © 2009, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

<sup>&</sup>lt;sup>1</sup>Operators can also 'ping' a phone to have it report back to a nearby tower, however this requires additional power from the phone and therefore typically is impractical for continuous location tracking.

2006). Generally this logging software records between one to four of the cellular towers with the highest signal strength, however, recent research suggests it is possible to obtain up to 2.5 meter accuracies if the number of detected towers is dramatically increased (Otsason et al. 2005).

**Human Movement Prediction** Dynamic Bayesian Networks (DBNs) have been widely used for quantifying and predicting human behavior. For analysis of human movement, typically these models involve location coordinates that are much more precise than cellular tower data, such as GPS data. These models are trained on general human movement (Ashbrook and Starner 2003) or more specific data such as transportation routes (Liao et al. 2004).

## Methods

## **Data Description**

Our data was generated from the phones of 215 subjects from a major urban city. After providing informed consent, these subjects were given phones that logged the ID of the four cellular towers with the strongest signal strength every 30 seconds. Additionally, the phones conducted Bluetooth scans every minute. Bluetooth beacons were deployed in the homes of each subject; as the beacons are detected only if the phone is within 10 meters of the beacon, detection implies the subject is at home. Additional data about the ambient audio environment was also collected, but not used for this analysis. The data was compressed on the mobile phone and sent via GPRS back to a central server after each day.

One unique difference with this data, compared to similar datasets reviewed above, is that every subject was randomly sampled from a particular city. By offering a smartphone and free service, over 80 percent of the randomly selected individuals agreed to participate in the study. The demographic information we have about the subjects is evenly distributed among ethnic groups and income levels, accurately reflecting the distribution that makes up the city's inhabitants. No longer constrained to the study of behavior of academics or researchers, our data represents an accurate depiction of the behavior of inhabitants in a major urban city.

#### Segmentation via Community Structure

Because each phone records the four towers with the strongest signal at 30 second intervals, this data can be represented as a cellular tower network (CTN) where each node is a unique cellular tower, an edge is placed if two towers ever co-occur in the same record, and each edge is annotated with the total amount of time (over all records) the pair co-occurred. A CTN is generated for each of the subjects, which includes every tower logged by the phone during the 5-month period. The nodes in the CTN that have the highest total edge weight (the node's "strength") correspond to the towers that are most often visible to the phone. Further, a group of nodes with a large amount of weight within the group, and less weight to other nodes, should correspond to a "location" where the user spends a significant amount of time. Figure 1 shows a 32-tower subgraph of one CTN, segmented into five such locations.



Figure 1: A 32-tower subgraph of one of our cellular tower networks, segmented into five "locations," clusters of nodes in which towers frequently coöccur in the phone's records.

To allow for a meaningful comparison, we use three qualitatively different heuristics for clustering nodes into locations.

**Ncut** The first segmentation algorithm depends on Shi and Malik's *normalized cut* (Ncut) criterion (Shi and Malik 2000), which, like many cut criteria, is NP-hard to optimize. Our implementation uses a spectral approach to find a bisection of the graph that minimizes the size of the normalized cut. Applied recursively, a graph can be split into a specified number of dense clusters. Although originally developed to segment images, the Ncut method can naturally be applied to networks.

**Q-Modularity** The second method, drawn from the large literature on detecting "communities" in complex networks (Newman 2006), depends on Newman and Girvan's popular *modularity* measure Q (Newman and Girvan 2004), which measures the density of clusters relative to a simple, randomized null model.

$$Q = \sum_{s=1}^{m} \left[ \frac{l_s}{L} - \left( \frac{d_s}{2L} \right)^2 \right] \tag{1}$$

where  $l_s$  is the number of edges between the nodes within cluster s, L is the total number of edges in the network, and  $d_s$  is the sum of degrees of the nodes in cluster s. While finding the segmentation that maximizes Q is NP-complete, there has been a significant amount of work towards this goal. Although also NP-hard to maximize, we use Clauset *et al.*'s greedy optimizer (Clauset, Newman, and Moore 2004), which has been shown to perform reasonably well on realworld data.

**Threshold Groups** The third method is a simple-minded heuristic: we first identify the nodes in the upper decile of "strength," and then perform a breadth-first search on the induced subgraph. Each connected component in this subgraph is labeled as a unique location, and all remaining nodes in the original graph placed in an additional group.

Although all based on somewhat similar principles, in practice these methods produce dramatically different segmentations of our CTNs. (This is in part because the first algorithm requires as input the number of segments to be found, unlike the other two.) One objective measure of these clusterings is to use independent information derived from the Bluetooth beacons.

#### **Inference via Bluetooth Beacons**

Bluetooth beacons have been installed in the homes of each subject in the study. Every minute the phone scans for visible Bluetooth devices and if a beacon is within 10 meters of the phone, it is logged as proximate. Creating training data from the set of cellular towers detected at the same time as the bluetooth beacons, we have used several methodologies to infer if a subject is at home given a particular set of visible cellular towers.

**Bayesian Posteriors** It is possible to calculate the posterior probability a subject is home,  $P(L_{home})$ , conditioned on the four towers currently detected by the phone,  $T_{abcd}$ , using the likelihood, the marginal and the prior probability of being at home (based on the beacon data).

$$P(L_{home}|T_{abcd}) = \frac{P(T_{abcd}|L_{home})P(L_{home})}{P(T_{abcd})} \quad (2)$$

Gaussian Processes While the naive Bayesian model above works well in many cases, simply using the ratio of tower counts co-present with the Bluetooth beacon tends to fail if the phone regularly moves beyond ten meters of the beacon while still staying inside the home. Instead of normalizing by total number of times each tower is detected, it is possible to obtain additional accuracy by incorporating the signal strengths from the detected towers. There are many models for signal strength of a single cellular tower, t.  $p_t(s_t|\mathbf{l})$ , one such model uses training data to estimate Gaussian distributions over functions modeling signal propagation from cellular towers (Schwaighofer et al. 2004). In our case, the training data comes from the signals of towers detected at the same time as the Bluetooth beacon in the subject's home, and the inference is binary (home or not home); however, these models are easily extendable for more broad localization.

**Deviations in Tower Signal Distributions** The two models above generate a probability of being at home associated with a single sample of detected towers (ie: the four tower IDs and their respective signal strengths). However, during the times when a subject is stationary, the phone continuously collects samples of the detected towers' signal strengths. These samples can form 'fingerprint' distributions of the expected signal strengths associated with that particular location. It is possible to detect deviations within these distributions of signal strengths using a pairwise analysis of variance (ANOVA) with the Bonferroni adjustment to correct for different sample sizes. Training the home distributions on the times when the beacon is visible (or if there are no beacons, on times when the subject is likely home such as 2-4am), an ANOVA comparing this home distribution with a distribution of recent tower signal strengths makes it possible to identify if the subject is truly at home, or is a nextdoor neighbor's house. In previous work, we have inferred such tower probability density distributions for office-level localization (Eagle and Pentland 2006).

### Prediction via Dynamic Bayesian Networks

The clusters of towers identified above can be incorporated as states of a dynamical model. Given a sequence of locations visited by a subject, we can learn patterns in their behaviour and calculate the probability of them moving to different future locations. We start with a baseline dynamical model and introduce additional observed and latent variables in order to model the situation more accurately.

The simplest dynamical Bayesian network we can use for location prediction is a Markov chain, in which the location  $y_t$  depends only on the location at the previous time step,  $y_{t-1}$ . The maximum likelihood transition probabilities  $p(y_t|y_{t-1})$  can easily be estimated. Given evidence that a user is in a particular location at time t, this allows us to calculate the  $\tau$ -step-ahead prediction  $p(y_{t+\tau}|y_t)$ .

We note that patterns of movement in practice are dependent on the time of day and the day of week. Subjects typically exhibit different dynamics on weekday mornings than on Saturday evenings, for example. Figure 2(a) shows an extended model where the probability of being in a location is also dependent on the hour of day  $h_t$  and the day of week  $d_t$ . In the experiments below, we code  $h_t$  to take on the values "morning", "afternoon", "evening" and "night", and code  $d_t$  to take on the values "weekday" or "weekend". After learning maximum likelihood parameters we can calculate the predicted density  $p(y_{t+\tau}|y_t, d_{t+1:t+\tau}, h_{t+1:t+\tau})$ for new observations from the same user.

# **X-Factors for Abnormality Modeling**

While there is strong structure in human behavior, there are also regular deviations from the standard routines. We incorporate an additional latent variable into our model to quantify the variation in behavior previously unaccounted for in the fully observed models above.

The model we use for this is shown in Figure 2(b). Here we factorize the location variable so that it depends on a hidden "abnormality" variable  $a_t$ . The model can now switch between "normal" and "abnormal" behaviour depending on whether  $a_t$  is 0 or 1 respectively.

We expect abnormal dynamics to be related to the normal dynamics but with a broader distribution. When estimating these dynamics, we therefore want to keep relevant structure in the dynamics (e.g. transitions between physically neighboring locations are still more likely), while allowing wider possibilities including non-zero probability of transitions not seen in the training data. We can achieve this effect by tying the parameters between the normal and abnormal transition probabilities such that  $p(y_t|y_{t-1}, d_t, h_t, a_t = 1)$  are a smoothed version of  $p(y_t|y_{t-1}, d_t, h_t, a_t = 0)$ . To smooth the transition matrices for every combination of  $d_t$  and  $h_t$  we add a small constant  $\xi$  to each entry in the matrix and renormalize.



Figure 2: Two DBN models used for location prediction. Shaded nodes are observed and unshaded nodes are latent;  $y_t$  denotes location,  $d_t$  denotes day of week,  $h_t$  denotes hour of day, and  $a_t$  denotes abnormal behaviour (all at time t). Panel (a) shows a fully observed model using contextual information, and panel (b) shows the X-factor model, where location is additionally conditioned on the latent abnormality variable.

Learning of this model can be done with expectationmaximization. We perform a standard E-step to calculate the probability of being in the normal or abnormal regime at each time frame, then modify the standard M-step to use the parameter tying above. In the experiments below, we set  $\xi = .1$  by hand, though in principle this parameter can also be learnt using EM. Increasing  $\xi$  effectively specifies that a sequence has to depart further from normal dynamics in order to be considered "abnormal".



Figure 3: Inferred points of abnormality using the X-Factor model. Each weekday the subject moves consistently between home (location 31) and work (location 15), but on the third day makes some extra, unusual journeys. The locations in this example were given by the local modularity segmentation method.

method	$\mu_{N_C}(\sigma)$	$P(C_{BT} \subset C_H)$	$\frac{N_{C_{BT}}}{N_{C_{H}}}$
Ncuts	20 (0)	.93	.0061
Q-Modularity	13.3 (11.7)	.86	.18
Threshold Groups	6.8 (13.7)	1.0	.045

Table 1: Segmentation Validation via Bluetooth Beacons.  $\mu_{N_C}$  is the average number of clusters generated by each segmentation method.  $P(C_{BT} \subset C_H)$  represents the probability that the set cellular towers associated the Bluetooth beacon at the subject's home,  $C_{BT}$ , is fully contained in a single cluster,  $C_H$ . The last column corresponds to the ratio of the actual number of home towers,  $N_{C_{BT}}$  to the number of home towers inferred by the different segmentation methods,  $N_{C_H}$ . A small number corresponds to incorporating a large number of towers within the home cluster.

# **Results & Discussion**

#### **Segmentation Validation**

We have shown how data collected from installed Bluetooth beacons can be used to create a known cluster of towers associated with each subject's home. We used this known cluster to validate each segmentation algorithm, selecting twenty locations for the Ncuts technique. Table 1 categorizes the community detection algorithms by how well they detected the "home" towers as defined by the Bluetooth beacons,  $C_{BT}$ . The home cluster of towers generated by the Threshold Groups technique incorporated  $C_{BT}$  for every subject,  $P(C_{BT} \subset C_H) = 100\%$ , while this was the case for the Q-Modularity technique only 86% of the time. However, the other important statistic is the ratio of the number of the Bluetooth home towers,  $N_{C_{BT}}$ , to the number of towers in the inferred home cluster,  $N_{C_H}^{D_1}$ . This ratio describes how many "extra" towers were included in the inferred home location; for example, the Q-Modularity home cluster has a ratio of .18, indicating that its home cluster contains approximately five times as many towers as needed. Despite averaging the most number of clusters, the Ncuts home cluster has a ratio of .0061, implying that a few large clusters tend to dominate these segmentations.

#### **Movement Prediction**

The three DBNs described above were trained on sequences of transitions between the locations that were inferred by each segmentation method. To compensate for the bias towards self-transitioning (at virtually every instance, the most likely event will be that the subject does not change locations), we compare the models success only on instances when a subject is about to transition between inferred locations. The DBNs are tasked with predicting the location where the subject is about to move. Table2 lists these prediction accuracies for the three segmentation methods and the two full-observed Markov models. While the X-factor model provides additional information about the the regularity a particular behavior, its accuracy is identical to the contextualized Markov model and was not included in the table. Of interest is that the highest accuracies did not come from the segmentation methods that provided the largest cluster

	Markov	Contextualized
method	Chain	Markov Chain
Ncuts	.932	.933
Q-Modularity	.953	.954
Threshold Groups	.992	.992

Table 2: Transition Accuracy. For every instance a subject moves between two clusters of towers, the DBN can be used to predict the subsequent cluster. The different accuracies between the segmentation methods are due to not only how well the clustering techniques performed at identifying the true salient locations, but also to the number and size of the clusters (described in Table 1). Given these high accuracies, contextual conditioning does not appear to provide significant improvement to the standard Markov model.

sizes (Ncuts), but rather with the smallest number of clusters (Threshold Groups). However, a direct comparison between these accuracies is not possible due to the differences in the dimensionality of the state spaces. A model with fewer inferred locations  $(N_C)$  should be expected to do better because it has less potential for a wrong prediction. In the extreme, a model with a single state will always be correct, yet obviously adds little value. Therefore, while the Threshold Groups segmentation method, with an average of 6.8 inferred salient locations ( $\sigma = 13.7$ ), generated accuracies of over 99%, future work in predicting location dwell times may provide more conclusive information about the dominance of one particular segmentation method over the others. Given the extremely high accuracies using an unconditioned Markov model, incorporating information about the time of day and day of the week unsurprisingly adds little additional value.

# **Future Work**

This paper has provided the groundwork for the design of increasingly sophisticated models that incorporate additional contextual and temporal variables and that can use demographic priors for bootstrapping. For example, if the discovered Bluetooth devices can be clustered based on copresense, it may be possible to classify particular Bluetooth phones as family, colleagues, and friends and incorporate the proximity of these individuals as observational variables. Additionally, the phones in this study also sample the ambient audio environmental periodically to detect the subjects' media consumption, information that should also make for an intriguing additional observed variable in the DBN. We are also planning on expanding these models to incorporate predicted dwell-times for locations. Lastly, we would like to explore the potential of using demographic bootstrapping to aid in efficient model parameterization as introduced in similar models (Liao et al. 2004).

We have demonstrated the potential to repurpose algorithms developed originally to quantify community structure within graphs to idenitfy salient locations within a cellular tower network. We have validated these unsupervised clustering algorithms on a known cluster of towers using the Bluetooth beacon installed in each of our randomly sampled



Figure 4: A sequence of transitions between clusters of towers corresponding to locations (top) and the average error rates for predicted transitions (bottom). The X-factor model was tested on approximately one month of movement segmented using Ncuts into 20 locations. While the top inferred location is 92% correct for this set of data, the subsequent location is in the top four locations over 99% of the time.

subjects' homes. The resultant set of inferred clusters of towers correspond to salient locations and are incorporated as states into our DBN models. We introduced the X-Factor model to detect behaviors that deviate from a given routine by incorporating an additional latent variable corresponding a normal / abnormal switch. We hope that these characterization and prediction methods for cellular tower data will contribute to the growing foundation that will support future pervasive computing applications.

### References

Ashbrook, D., and Starner, T. 2003. Using gps to learn significant locations and predict movement across multiple users. *Personal and Ubiquitous Computing* 7:275–286.

Chen, M.; Sohn, T.; Chmelev, D.; Haehnel, D.; Hightower, J.; Hughes, J.; LaMarca, A.; Potter, F.; Smith, I.; and Varshavsky, A. 2006. Practical metropolitan-scale positioning for gsm phones. *Ubicomp 2006, Lecture Notes in Computer Science* 4206:225–242.

Clauset, A.; Newman, M. E. J.; and Moore, C. 2004. Finding community structure in very large networks. *Physical Review E* 70(6).

Davis, M.; King, S.; Good, N.; and Sarvas, R. 2004. From context to content: leveraging context to infer media metadata. *Proceedings of the 12th annual ACM international conference on Multimedia, October 10-16, 2004, New York, NY, USA.* 

Eagle, N., and Pentland, A. 2006. Reality mining: sensing complex social systems. *Personal and Ubiquitous Computing* 10:255–268.

González, M. C.; Hidalgo, C. A.; and Barabási, A.-L. 2008. Understanding individual human mobility patterns. *Nature* 453(7196):779–782.

LaMarca, A.; Chawathe, Y.; Consolvo, S.; and Hightower, J. 2005. Place lab: Device positioning using radio beacons in the wild. *Pervasive 2005, LNCS* 3468:116–133.

Liao, L.; Patterson, D.; Fox, D.; and Kautz, H. 2004. Learning and inferring transportation routines. *Proceedings of the Nineteenth National Conference on Artificial Intelligence* 348–353.

Newman, M., and Girvan, M. 2004. Finding and evaluating community structure in networks. *Physical Review E* 69.

Newman, M. 2006. Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*.

Onnela, J.; Saramaki, J.; Hyvonen, J.; Szabo, G.; Lazer, D.; Kaski, K.; Kertesz, J.; and Barabasi, A.-L. 2007. Structure and tie strengths in mobile communication networks. *Proceedings of the National Academy of Sciences* 104:7332–7336.

Otsason, V.; Varshavsky, A.; LaMarca, A.; and de Lara, E. 2005. Accurate gsm indoor localization. *Ubicomp 2005*, *LNCS* 141–158.

Schwaighofer, A.; Grigoras, M.; Tresp, V.; and Hoffmann, C. 2004. Gpps: A gaussian process positioning system for cellular networks. *Advances in Neural Information Processing Systems* 16.

Shi, J., and Malik, J. 2000. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Learning* 22(8):888–905.

Sohn, T.; Varshavsky, A.; LaMarca, A.; and Chen, M. 2006. Mobility detection using everyday gsm traces. *Ubicomp 2006, Lecture Notes in Computer Science* 212–224.