# Predicting your footsteps on Facebook

**Manav Verma**
mverma4@ncsu.edu

**Sameer Sharma**
ssharm20@ncsu.edu

**Shivam Gulati**
sgulati2@ncsu.edu

## Abstract

Facebook check-in activity has become a source of recommendations and targeted digital advertising for the customer and businesses these days. This paper aims to improve upon the check-in prediction accuracy on the dataset provided by Facebook. In this paper, we use different machine learning classifiers and compare their results to find the optimal model.

## 1 Background and Introduction

### 1.1 Introduction

As a person, what if you knew where you could be heading to next? As a business or a place of interest, how about you had an idea how many customers or people could be interested for the same. We travel, eat out at the restaurants, visit other places of interest, and the generation being Facebook friendly, check-in at different places.

A system which could classify and analyze social media check-in data to find patterns on how users check-in activity would be beneficial for the customers and businesses. Although, achieving this would never have a hundred percent accuracy, but a good prediction would certainly do wonders for both. Our idea is to predict where a user could be checking into from their location and timestamp. This project was a coding challenge on kaggle.com[1], and the dataset we are using is from the website itself.

### 1.2 Applications and business value

Classification and prediction becomes the basis of recommendation of places of interests. Such systems have already made a major impact in online service industry. This would provide great business values to the businesses from digital advertising point and view as well as to the customers. Few of the applications are highlighted below:

- See the most common places of interests as per the geographical area.
- Allow businesses to promote themselves as per customer check-ins and find ways to improve and work on customer retention
- Realize how customer visit places as per the seasons, days, hours, etc. More check-ins in evening could be a sign of restaurant outings and allow to find favorable eating joints.
- Allow users to find the most visited hangout places around.
- Allow tourists to better plan their trips highlighting how people go about visiting a city and area as per different time of the day.
- Customized offers and personalized customer based advertising on social media and other networks.

### 1.3 Related work

This has been a booming area, and not only specifically Facebook, but location prediction in general has been widely researched. We did find some interesting related work which was a guide to us working on this project.

The paper "A Random Walk Around the City: New Venue Recommendation in Location-Based Social Networks"[2] talks about how large datasets from apps and devices are extracted and analyzed with algorithms and how efficient it proves for location recommendation. The authors have focused on the datasets of Foursquare and Gowalla for this one.

An interesting read was the mobile data mining framework NextLocation[3], which predicts the next place, taking privacy and many other factors into our account. It talks about how the data was preprocessed and transformed with sequences and timestamps for the best prediction and analysis.

With data sharing, also comes the security aspect. Users do not prefer to share all the personal information without any consent. The paper "Predicting Users' Motivations behind Location Check-Ins and Utility Implications of Privacy Protection Mechanisms"[4] talked about the privacy concerns and suggested how generalization of check-ins by finding the motivation behind it solves the purpose and yet, helps us to analyze the way we want it. In our dataset, we only had the information required for the project and the names were generalized as ids beforehand.

## 1.4    The Dataset

The dataset[6] we are using was provided with the project on kaggle.com. We were provided with both the testing and training data. The data files contain the following columns:

```
row_id: id of the checkin event
x: x coordinate
y: y coordinate
accuracy: accuracy of the location
time: timestamp of the checkin
place_id: id of the business
```

The place-id is what we are trying to predict. It could be a most favorable place_id or even extend to a numbered ranking list as per different patterns. One of the major challenges was to make the data more meaningful by finding patterns on features. Hence, data preprocessing and feature extraction is a critical step in this scenario.

## 2    Method

For this task of predicting check-ins of the user, Facebook has provided more than 30 million previous check-ins in a 10 by 10km grid. As described above the dataset consists of only two csvs one for training and other to test our model. The train csv consists of features: x, y, accuracy, timestamp and place_id. Our task is to predict top three probable place_ids or businesses a user checks into based on spatial and temporal information given in testing data. This is what we are trying to achieve in nutshell in this project.

To get started with the project and explore further the data in hand we set out by applying some of the most common algorithms like Linear Regression to the train data. These algorithms resulted in very bad results as we have around 100K place_ids in our dataset. Most of the supervised learning algorithm fails to work successfully with that number of classes. Thus, we needed to move on a different path or transform our data in a way these algorithms can be applied to it.

### 2.1    Underlying challenges in dataset

The sheer number of place_ids or classes to predict presents the biggest challenge in this project. Facebook or any other social media site that presents users with functionality like check-in do not put any restrictions on places a user can check into. This results in many places that can be checked into and hence the large number of place_ids present in our data is a common problem faced during the prediction of check-ins for a user.

The other challenging aspect of the given dataset is that definition of accuracy and timestamp column is left open to the interpretation of the candidate. This is done on purpose just to make the

task more challenging. Among these two attributes timestamp is easier to interpret and to model how it can affect check-ins of a user. We know that for some businesses like restaurants lunch and dinner hours are busy hours of the day whereas if someone is checking into a mall it can be at any time of a day. Same way some places witness much more crowd during the weekend as compared to weekdays. The second attribute accuracy can be the accuracy with which x and y coordinate values of the check in is recorded. Though, exact distribution function for accuracy and its correct usage in the model is difficult to access.

## 2.2    Feature extraction

In the data set given to us we have five columns. The first four columns have been used to train the model to predict the place_id. We extracted some additional features from the given data that are discussed below:

- The first feature is 'x': This is x coordinate of the check-in, this feature is important in the data set in predicting the place_id, combined with the y coordinate it is used to cluster the check-in position. No data cleaning was required for this feature.
- The second feature is 'y': Similar as first feature, it is the y coordinate of the check-in.
- The third feature is 'accuracy'. This represents the accuracy of the x and y coordinates. The description for this was vague, thus we decided to use it directly in the random forest model and skipped it in the KNN model. No data cleaning was required for this feature.
- The fourth feature is 'hour': This is extracted from the timestamp and denotes the hour of the 0-24 at time of check-in.
- The fifth feature is 'weekday': This denotes the day of the week 1-7 for check-in.
- The sixth feature is 'month' of the year in range 1-12.
- The seventh feature we used is 'year'.
- The eighth feature is 'day': This is the day of the year (1-365).
- The ninth feature is 'place_id': This is what we are trying to predict in our model.

## 2.3    Exploratory Analysis

Due to very large number of check-ins (> 30 million) provided to us in the training data, it very difficult to visualize data for whole grid and look for patterns in it. Hence, we can take data in a smaller grid and try finding patterns or clusters in it. We read numerous approaches taken to solve this problem and the first recommended step was to partition the grid or map in smaller unit grids. This partitioning makes sense in our case because check-in for a given place is generally done nearby that location and does not vary by large distances. Thus, given the position of a person in terms of x and y coordinates, we can build an imaginary smaller grid around that x, y location and apply the classification or clustering techniques to only data points in this grid.

To reinforce our understanding of the problem and to visualize the data we start by plotting all the check-ins within a smaller grid of 500 X 500 meters taken at random from the given larger grid. We have plotted only place_ids that have more than 100 check-ins to visualize the clusters.

Some clusters are evident in the plot shown below (Figure 1) but most of these clusters overlap with each other and difficult to separate. To make these clusters separable and more evident we tried using one more feature "hour of the day" as the third dimension for our plot (Figure 2).

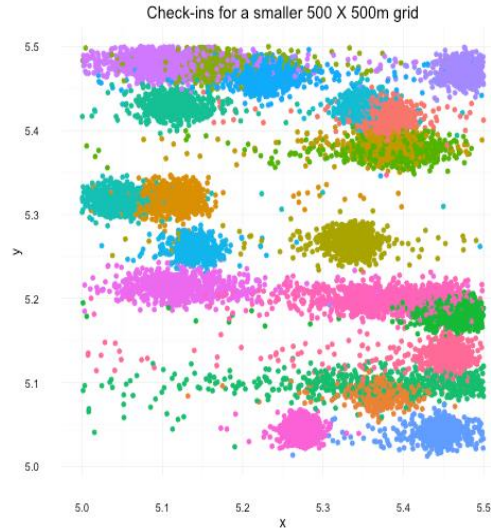Check-ins for a smaller 500 X 500m grid



Figure 1: Place_id clustered on location

Addition of third dimension helps and we can see that our assumption that hour of day affects the check-ins for a place is valid. We tried plotting the same data using "weekday" feature that resulted in similar plot. These plots confirm our understanding that check-ins from user depends on the different time components like hour of the day and weekday.
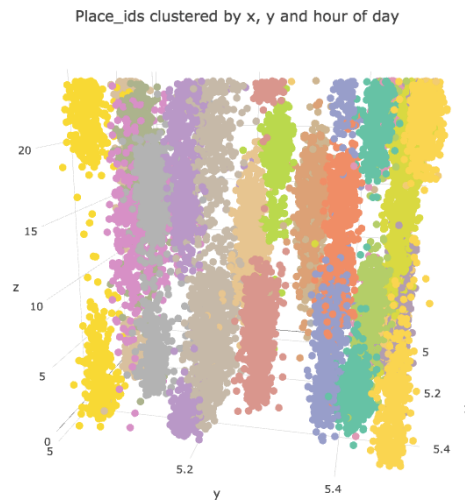
Place_ids clustered by x, y and hour of day



Figure 2: Place_id clustered on location and hour of day

## 2.4    Algorithms and Data classification techniques

For classification task, we deal with one row of the testing data at a time. As explained in the last section it is safe to assume that a check-in for a place will be done within some specified distance from the location of that place. Thus, the first step we do while processing a row of test data is draw an imaginary grid of 250 X 250 meters around that x and y coordinates. Once we have this grid in place we fetch all the data related to the points inside this grid from our training data. Now for doing the classification we will consider only these rows because the check-in location that a user can select lies in this grid. For classification, we can use multiple classifiers and techniques like clustering with KNN. As our initial models, we chose KNN and Random Forest for the same.

### 2.4.1    K Nearest Neighbors

Once we have a smaller grid of 250 X 250 meters in place first thing that comes to mind is KNN for the classification task. It is very easy to implement and give good results. The only tricky part

about applying KNN is finding out the optimal weights for the variables used. We have used hit and trial method to optimize our model. We plan to use data exploratory techniques in future to narrow down on optimal weights for KNN in final version of the report.

### 2.4.2 Random Forest

Random Forest was our second choice for the classifier as it is efficient and generally results in more accurate results. The performance factor of random forest is important for us as we are doing classification task on the fly. We chose random forest also because it gives us estimate of the importance of different variables in classification task this helped us in fine-tuning our model for other classifiers as well. We tried different flavors of Random Forest available in Python and achieved best results using sklearn random forest classifier.

### 2.4.3 Boosted Trees

To improve the accuracy further, we tried boosted trees i.e. tree ensemble model for classification and regression trees (CART). We used XGBoost library[7], short for "Extreme Gradient Boosting", where the term "Gradient Boosting" is proposed in the paper Greedy Function Approximation: A Gradient Boosting Machine, by Friedman. We tried boosted trees as our third classifier as it was mentioned in the blogs and discussions on this competition to provide good results.

## 3 Experimentation and Results

We started with three above models and the accuracy we obtained for the initial test run were 0.40, 0.48 and 0.46 for K Nearest Neighbors, Random Forest and Boosted Trees respectively. These are just initial accuracy values we obtained for the models. We further plan to improve the accuracy by tweaking the parameters and using ensemble models. We also plan to evaluate the results we obtained and access the advantages and disadvantages of each model.

## 4 Conclusion and Future Work

We obtained good results using very simplistic model and strategy but we are still lagging in terms of accuracy, as the winner of the competition achieved accuracy of 0.58. We plan to improve on the existing models and try several other classifiers to improve our result. We would also include the elaborate evaluation of the work done and results obtained to get a better understanding of the results.

## 5 References

[1] Kaggle project challenge: https://www.kaggle.com/c/facebook-v-predicting-check-ins
[2] A Random Walk Around the City: New Venue Recommendation in Location-Based Social Networks: Anastasios Noulas, Salvatore Scellato, Neal Lathia, Cecilia Mascolo
[3] Where will you go? Mobile Data Mining for Next Place Prediction: Joao Bartolo Gomes, Clifton Phua, Shonali Krishnaswamy
[4] Predicting Users' Motivations behind Location Check-Ins and Utility Implications of Privacy Protection Mechanisms: Igor Bilogrevic, Kevin Huguenin, Stefan Mihaila, Reza Shokri, Jean Pierre Hubaux
[5] A data mining approach for location prediction in mobile environments: Gokhan Yavas, Dimitrios Katsaros, Ozgur Ulusoy, Yannis Manolopoulos
[6] Project dataset: https://www.kaggle.com/c/facebook-v-predicting-check-ins/data
[7] XGBoost Library: http://xgboost.readthedocs.io/en/latest/model.html

## A. Appendix

As per our initial plan, we have completed data preprocessing and build models on top of it. For the final phase, we plan to improve the accuracy of these models and evaluate the results obtained. All the team members contributed equally to the planning and implementation phases.