

# New York City’s Citi Bike System Simulation

Saul Toscano-Palmerin

Cornell University

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Simulation</b>	<b>2</b>
2.1	Run the Simulation in Python . . . . .	2
<b>3</b>	<b>Solution of the Citi Bike Problem</b>	<b>4</b>
<b>4</b>	<b>Run SBO, KG, and EI on the Citi Bike Problem</b>	<b>5</b>

# 1 Introduction

We consider a realistic problem, using a queuing simulation based on New York City's Citi Bike system, in which system users may remove an available bike from a station at one location within the city, and ride it to a station with an available dock in some other location within the city. The optimization problem that we consider is the allocation of a constrained number of bikes (6000) to available docks within the city at the start of rush hour, so as to minimize, in simulation, the expected number of potential trips in which the rider could not find an available bike at their preferred origination station, or could not find an available dock at their preferred destination station. We call such trips "negatively affected trips."

## 2 Simulation

We simulated in Python the demand of bike trips of a New York City's Bike System on any day from January 1st to December 31st between 7:00am and 11:00am. We used 329 actual bike stations, locations, and numbers of docks from the Citi Bike system, and estimated demand and average time for trips for every day in a year using publicly available data of the year 2014 from Citi Bike's website Citi (2015).

We simulate the demand for trips between each pair of bike stations on a day using an independent Poisson process, and trip times between pairs of stations follows an exponential distribution. If a potential trip's origination station has no available bikes, then that trip does not occur, and we increment our count of negatively affected trips. If a trip does occur, and its preferred destination station does not have an available dock, then we also increment our count of negatively affected trips, and the bike is returned to the closest bike station with available docks.

We divided the bike stations in 4 groups using k-nearest neighbors, and let  $x$  be the number of bikes in each group at 7:00 AM. We suppose that bikes are allocated uniformly among stations within a single group.

The simulation has five steps: (i) randomly choose a day  $i$  from January 1st to December 31st; (ii) simulate the total demand denoted by  $w$  as a Poisson random variable with mean  $\lambda_i$ , where  $\lambda_i$  is the total number of bike trips of day  $i$ ; (iii) simulate the demands of the directed pairs of bike stations  $(j, k)$  as a multinomial random variable with parameters  $(w, (p_{j,k} : 1 \leq j, k \leq 329))$  where  $p_{j,k}$  is the proportion of bike trips from  $j$  to  $k$  on that day; (iv) simulate the arrival times to the bike stations as uniform random variables from 7:00am to 11:00am; (v) if a rider does not find an available bike, then he will leave the system, and we increment our count of negatively affected trips. If a rider does find an available bike, and his destination station does not have an available dock, then we also increment our count of negatively affected trips, and he will return the bike to the closest bike station with available docks.

### 2.1 Run the Simulation in Python

The Python file of the simulation is `simulationPoissonProcessNonHomogeneous.py`, and the function `negativelyAffectedTrips` computes the number of negatively affected trips.

---

```
def negativelyAffectedTrips (T,N,X,m,cluster,bikeData,parLambda,nDays,A,
                             poissonArray,timesArray,ind=None,randomSeed=None,
                             nStations=329):
```

```

"""
Counts the number of negatively affected trips.
We divide the bike stations in  $m$  groups according to K-means algorithm.
The bikes are distributed uniformly in each group.

Args:
    T (int): Duration of the simulation in hours
              (it always starts at 7:00am).
    N (numpy array): Vector  $N(T, A_{\{i\}})$ .
    X (numpy array): Vector with the initial configuration of the bikes.
    m (int): Number of groups formed with the bike stations.
    cluster (List[List[List[int, float, float]]]):
                Contains the clusters of the bike stations
                with their ids, and geographic coordinates.
    bikeData (numpy array): Matrix with the ID, numberDocks, Latitude,
                            and longitude of each bike station.
    parLambda (numpy array) : Vector with the parameters of the
                              Poisson processes  $N(T, A_{\{i\}})$ .
    nDays: Number of different days considered in the simulation
           (i.e. 365).
    A (List[List[Tuple(int, int)]]):
        List with subsets of pair of bike stations.
    lamb (List[numpy array]): List with the parameters of the Poisson
                              processes  $N(T, (i, j))$ ,
                              (i.e.  $j$ th entry of lamb is the Poisson process with parameter
                               $\text{lamb}[0][0, j]$  between stations  $\text{lamb}[0][1, j]$  and  $\text{lamb}[0][2, j]$ ).
    poissonArray (List[List[numpy array]]):
        List with the parameters of the Poisson processes  $N(T, (i, j))$ ,
        where the  $j$ th entry of poissonArray are the parameters of the
        Poissons processes of day  $j$  between each pair of bike stations.
        (i.e., the parameter of the Poisson process on day  $j$  between
        stations  $\text{lamb}[j][0][1, 1]$  and  $\text{lamb}[j][0][2, 1]$  is  $\text{lamb}[j][0][0, 1]$ .
        This is a sparse representation of the original matrix, and so
        if a pair of stations doesn't appear in the last list, its PP
        has parameter zero.
    timesArray (List[List[numpy array]]):
        Similar tan poissonArray, but with the mean times of traveling
        between the stations.
    ind (int or None): Day of the year when the simulation is run.
    randomSeed (int): Random seed.
    nStations (int): Number of bike stations.

Returns:
    int: Overall number of negatively affected tripes multiplied by  $-1$ .
"""

```

---

For our experiments, we created the data with  $m = 4$  and  $nDays = 365$ . The repository already includes the files for poissonArray, timesArray, lamb, A, bikeData, and cluster. We can create these parameters in the following

way.

---

```
fil="poissonDays.txt"
fil=os.path.join("NonHomogeneousPP2",fil)
poissonParameters=np.loadtxt(fil)

poissonArray=[[[] for i in xrange(nDays)]
timesArray=[[[] for i in xrange(nDays)]

for i in xrange(nDays):
    fil="daySparse"+"%d"%i+"ExponentialTimesNonHom.txt"
    fil2=os.path.join("SparseNonHomogeneousPP2",fil)
    poissonArray[i].append(np.loadtxt(fil2))

    fil="daySparse"+"%d"%i+"PoissonParametersNonHom.txt"
    fil2=os.path.join("SparseNonHomogeneousPP2",fil)
    timesArray[i].append(np.loadtxt(fil2))

TimeHours=4.0
poissonParameters*=TimeHours

Avertices=[[[]]
for j in range(numberStations):
    for k in range(numberStations):
        Avertices[0].append((j,k))

f = open(str(4)+"-cluster.txt", 'r')
cluster=eval(f.read())
f.close()

bikeData=np.loadtxt("bikesStationsOrdinalIDnumberDocks.txt", skiprows=1)
```

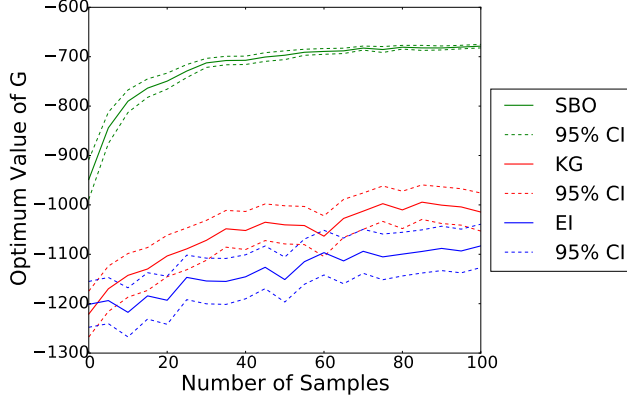
---

### 3 Solution of the Citi Bike Problem

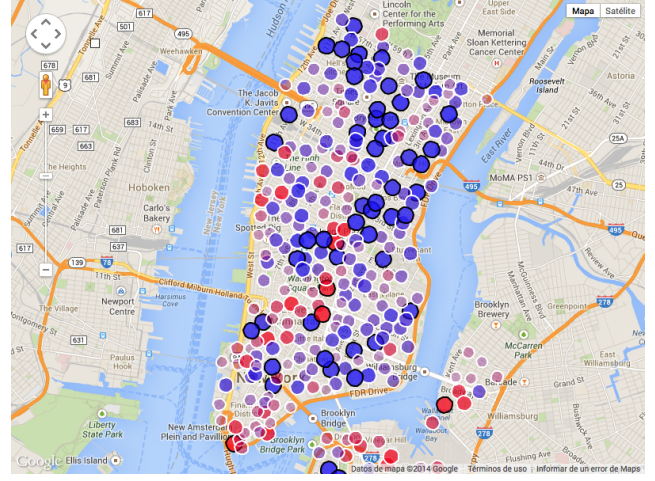
We solve the Citi Bike problem, arising in the design of the New York City's Citi Bike system, using three Bayesian optimization algorithms. We use Stratified Bayesian Optimization of (Toscano-Palmerin and Frazier, 2016), the Knowledge-Gradient policy of Frazier *et al.* (2009) and Expected Improvement criterion (Jones *et al.*, 1998).

We define our objective function by  $G(x) = \mathbb{E}[f(x)]$  where  $-f$  is the number of negatively affected trips between 7:00am to 11:00am, and  $x$  is the number of bikes in each group of bike stations at 7:00 AM.

Figure 1a compares the performance of Stratified Bayesian Optimization (SBO), Knowledge Gradient (KG) and Expected Improvement (EI), plotting the number of samples beyond the first stage on the  $x$  axis, and the average true quality of the solutions provided,  $G(\arg\max_x \mathbb{E}_n[G(x)])$ , averaging over 300 independent runs of the three algorithms. We see that SBO was able to quickly find an allocation of bikes to groups that attains a small expected number of negatively affected trips.



(a) Performance comparison between SBO and two Bayesian optimization benchmark, the KG and EI methods, on the Citi Bike Problem



(b) Location of bike stations (circles) in New York City, where size and color represent the ratio of available bikes to available docks.

Figure 1: Performance results for the Citi Bike problem (plot 1), and a screenshot from our simulation of the Citi Bike problem (plot b).

## 4 Run SBO, KG, and EI on the Citi Bike Problem

The scripts are `SBOcitiBikeExample.py`, `KGcitiBikeExample.py`, and `EIcitiBikeExample.py`. The arguments of these scripts are:

- (int): Random seed.
- (int): Number of points in the training data.
- (int): Number of samples to estimate  $F$  when SBO is used, and  $G$  otherwise.
- (int): Number of iterations of the algorithm.
- (char): T if the code is run in parallel when we use multiple starting points when optimizing VOI and  $a_n$ ; F otherwise.
- (int): Number of the starting points used to optimize the VOI and  $a_n$ , which is given only when the previous argument is T.

We can run those procedures during 15 iterations with 30 training data points, 10 samples to estimate  $F$  or  $G$ , 1 as the random seed, and using 10 different starting points when optimizing VOI and  $a_n$ , in the following way:

by writing:

- `python SBOcitiBikeExample.py 1 30 10 15 T 10`
- `python KGcitiBikeExample.py 1 30 10 15 T 10`
- `python EIcitiBikeExample.py 1 30 10 15 T 10`

The output will be saved in the directory

FinalNonHomogeneous011116Results15AveragingSamples30TrainingPoints/SBO/1run.

Please refer to <https://github.com/toscanosaul/BGO/blob/master/BGO.pdf> for a complete description of the output.

## References

- Citi (2015). Citi bike website. <https://www.citibikenyc.com/>, accessed November 2015.
- Frazier, P., Powell, W., and Dayanik, S. (2009). The knowledge-gradient policy for correlated normal beliefs. *INFORMS journal on Computing*, **21**(4), 599–613.
- Jones, D. R., Schonlau, M., and Welch, W. J. (1998). Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, **13**(4), 455–492.
- Toscano-Palmerin, S. and Frazier, P. (2016). Stratified Bayesian Optimization. *arxiv 1602.02338*.