

**FOOD ESTABLISHMENT INSPECTION FROM ANALYSE BOSTON****PAPER ON MY MEASURES**Introduction:

The Health Division of the Department of Inspectional Services ensures that all food establishments in the City of Boston meet relevant sanitary codes and standards. Businesses that serve food are inspected at least once a year, and follow-up inspections are performed on high-risk establishments. The dataset being referred to in this paper is a legacy dataset containing records of individual inspections and results, with a total of 558545 cases and 30 columns/variables related to business names, locations, sanitary violation levels, reasons of violation, comments of improvement, etc., for those food establishments in Greater Boston Area.

Latent constructs or measures are theoretical in nature; they cannot be observed directly and, therefore, cannot be measured directly either. To measure a latent construct, researchers capture indicators that represent the underlying construct. I was intrigued by the possible measures that can be derived from the variables of the aforesaid dataset, so after much analysing I decided on exploring how the health department assesses the violation of a given business chain.

I trust we can allude to the violation description given for each dimension of infringement and the remarks expressed by the health department. I likewise chose investigating the Violation levels and statuses as for the permit statuses of each business chain with the goal that I can draw a superior image of the Greater Boston region and can remark about an area's strength.

I found this measure intriguing on the grounds that my entire dataset depends on investigation of sustenance foundations and I need to realise how does the Health division of Boston concocts their violation level, what is that factor which gives these eatery's their stars or takes them away.

I want to assess the degree of every violation, I want to see what violation houses what kinds of comments or description. This gives us the possibility of the area also. The area with a particular arrangement of infringement can enable us to comprehend the format of its neighbourhood. For instance, if every one of the eateries in South Boston territory has plumbing issues, one can comprehend that this zone has poor pipes or this current zone's proprietors don't generally deal with their plumbing issues.

On the off chance that a specific neighbourhood has the most elevated number of gravely damaged eateries we can express that the zone is ill-advised or has terrible natural surroundings. This makes us layout a picture of the all inclusive community around us.

#### Methods:

##### ***Datasets and Sources:***

Our spatial unit of investigation is the Violation measure in various eateries all through the Greater Boston region. Data on Food Establishment Inspection was made accessible by the US Health Department on Analyse Boston-in 2016 chain - violation level inside a metropolitan territory (here Boston) in a specific timeframe. Information to delineate city through shapefiles was obtained by me from Harvard's Dataverse site under the file name Tracts BARI dataset - 2010 chain.

##### ***Preliminary Analysis and Conversions:***

I attempted to utilise the license numbers (LICENSE NO) as an id regarding its status (LIC STATUS). I endeavoured to perceive what number of licenses are active and non-active this gave me an insight into structuring the latent construct. Utilising the total capacity determined the tally of the two kinds of licenses. Utilising the aggregate function I endeavoured to delineate territory's

postal division (ZIP) to its property id with the goal that we can put every business chain region shrewd. The business name (businessName) variable gives us the check of the specific property id (PROPERTY\_ID) with the absolute number of eateries under that business chain. I also converted the license status into numbers and calculated the “sum” and “mean” for both the cases to check how many percentage of business chains had active licenses and how many didn’t.

### ***String Analysis with Bigrams:***

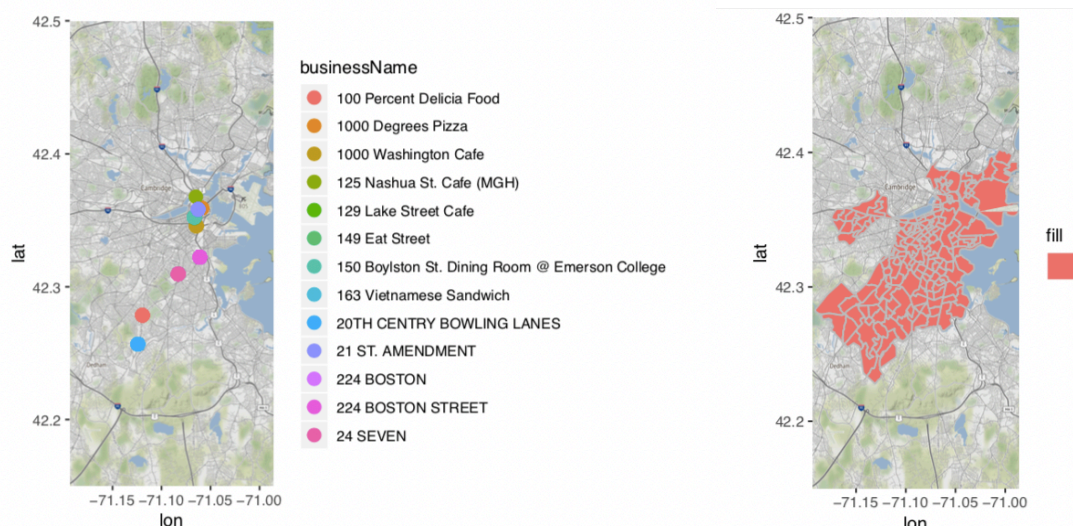
I utilised the Violation description (ViolDesc) and the remarks (Comments) given by the Health division to set up the reason concerning why a specific eatery gets its infringement and furthermore endeavouring to outline a particular region of the city to roughly implement our final latent construct. I converted the description into factors. Imported all of the comments and violation description into a tidy text format and converted them into tibbles for easy assessment. I converted the column of comments into factors and separated the food, management and pest issues related complaints and did the same for Violation Description.

Reading txt document and converting into tibbles was finished by changing over all the depiction and remarks to txt record and after that bringing in them as four separate records. Bigram analysis was performed to develop the dormant build. They were put away under the accompanying titles; clean\_issue, management\_issue, food\_issue, animal\_issue. I for the most part determined the "count" for every one of these keywords to see which was the primary driver of violation and after that mapped it to the latitude and the longitude in a like manner for a specific neighbourhood. I attempted to delineate these issues with every one of the areas in Boston. (Latitude and Longitude)

.

***Mapping measure using shape-files:***

I imported the BARI tracts dataset and stacked the shape files and assessed it regarding my measure. The variable utilised in BARI was CT\_ID\_10. It geocoded all the business chains with their latitudes and longitudes and subsequently delineated my construct on the Boston map. I plotted the primary inactive build which incorporated the property id, (Property\_id) the postal district, (ZIP) the license status (LIC-STATUS) and the violation level (Viol Level) concerning all the area zones in Boston. (CITY) . At that point the second half of the measure was plotted disclosing to me which regions in boston has cleanliness related issues and which business chains have pest related issues because through my bigram analysis earlier the pest and food related violation seemed to be the most prevalent issues for violations.

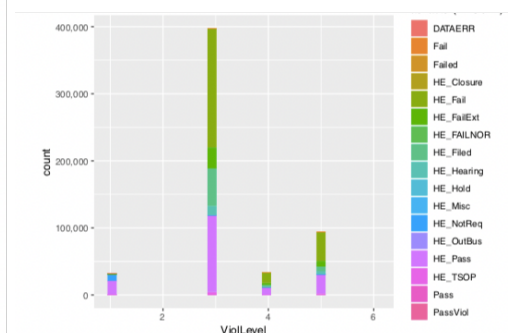
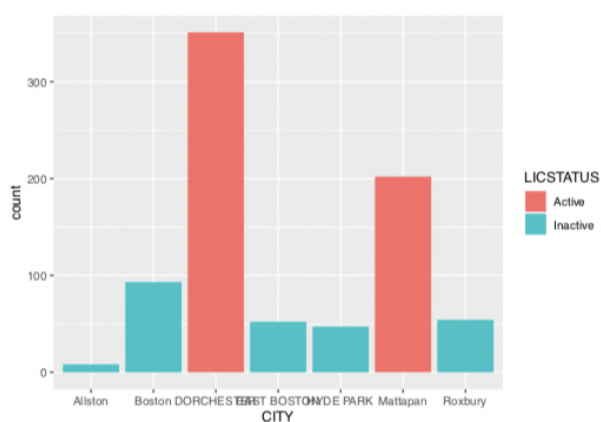


***Inferential statistics:***

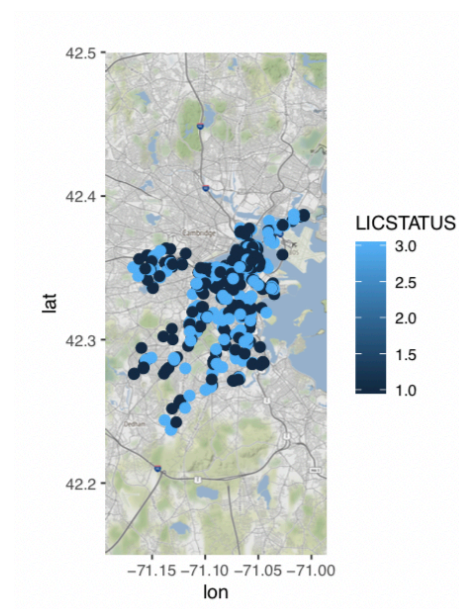
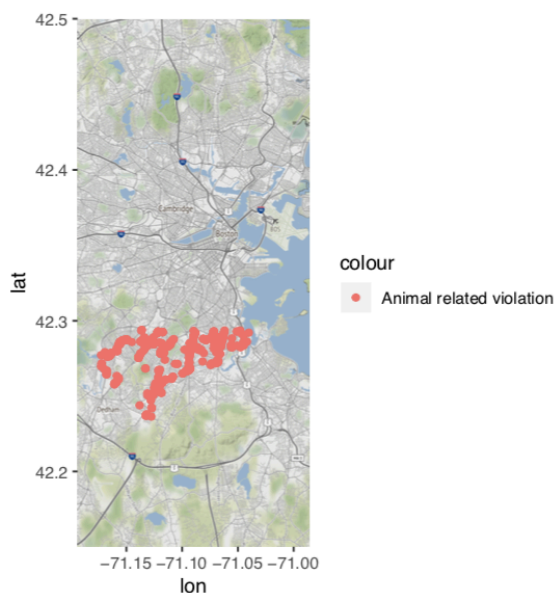
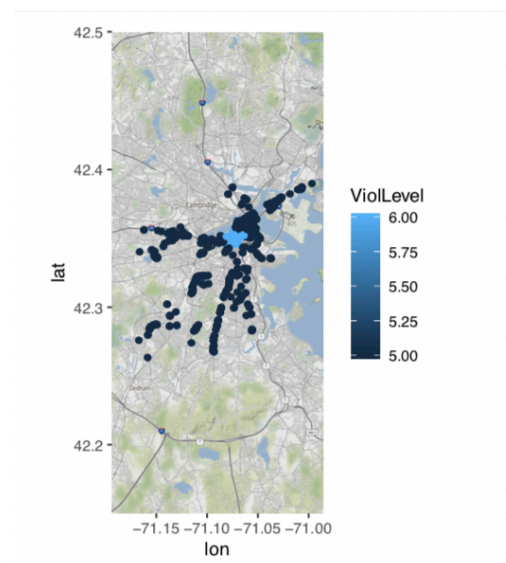
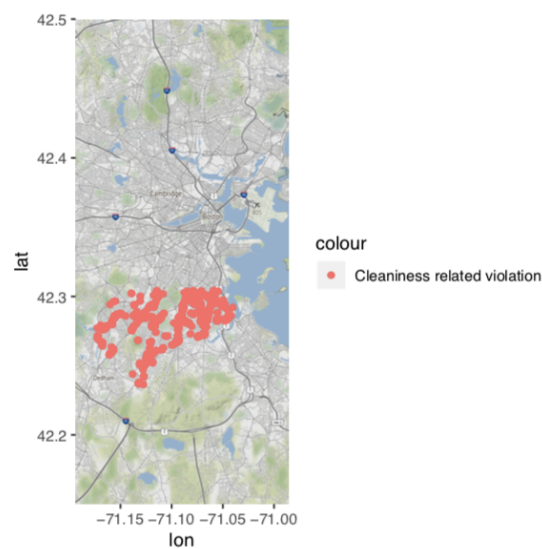
To upgrade my measures, I utilised inferential insights and statistics , t-test and ANOVA test to convey more profundity to my measure with the goal that I can portray my examination and exhibit the appropriation of basic measures over the city.

**Results:**

A vast segment of these business chains have Violation Level 3 and Level 1 when changed over to numeric from character structure, there is a range from Level 1 to Level 5 which helped me classify the infringement need and prioritise violation levels. The tally of active and non-active licenses was astounding as was the violation level being housed by these licenses since they uncovered a non-reliable and least anticipated outcome, which can be utilised by the city council to investigate the Health department and it's proceedings. Total number of non-active licenses are 3289, which means that 47 % of the total licenses given out to the business chains are inactive. Total number of active licenses are 3605, which means that 52 % of the total licenses given out to the business chains are active.



The analysis of bigrams and the word count with respect to the business chains displayed the common issue of violation at these places and also gave me the count of these eateries so that one can map them and focus on making policies for the betterment of these business chain.



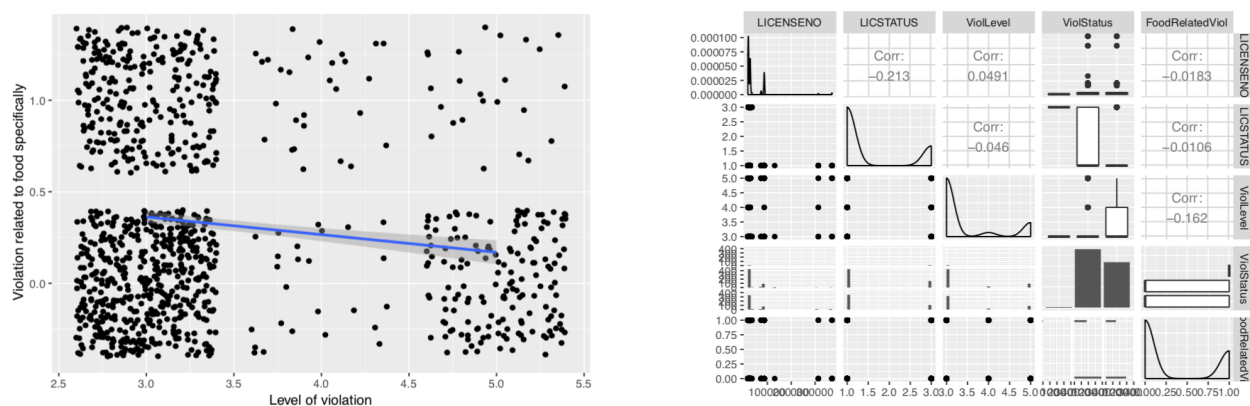


In the wake of experiencing bigrams I understood that the ill-advised upkeep of walls and roofs have the most noteworthy tally. Next came floors and storage, followed by improper storage and use of utensils and different machines. The cleanliness issue has unclean floor as the most elevated check and it appears as one of the fundamental issues of infringement for an eatery. Followed by unclean interiors and exteriors. The mouse droppings has the highest count in the animal issue wing, surprisingly there is a separate count for rodent droppings. These are trailed by lack of pest control. I observed that the food establishments that have cleanliness related issue also definitely have animal and pest related issue in common, we can see that through the graphs. As we can see from the map above that the focal boston has the most damaged natural pecking orders contrasted with different zones.

## # A tibble: 15 x 3				## # A tibble: 15 x 3				## # A tibble: 15 x 3			
##	word1	word2	n	##	word1	word2	n	##	word1	word2	n
##	<fct>	<fct>	<int>	##	<fct>	<fct>	<int>	##	<fct>	<fct>	<int>
##	1 improper	maintenance	867	##	1 clean	floor	134	##	1 contact	surfaces	650
##	2 walls	ceilings	460	##	2 clean	interior	93	##	2 food	contact	650
##	3 ceilings	improper	214	##	3 interior	exterior	53	##	3 surfaces	clean	367
##	4 floors	improper	197	##	4 properly	clean	53	##	4 food	protection	196
##	5 improperly	maintained	85	##	5 ice	machine	45	##	5 clean	food	98
##	6 storage	improperly	85	##	6 clean	organize	37	##	6 container	labels	84
##	7 maintained	improper	43	##	7 clean	sanitize	37	##	7 food	container	84
##	8 improper	cleaning	28	##	8 cooking	equipment	35	##	8 surfaces	design	84
##	9 improper	storage	20	##	9 hot	line	35	##	9 protection	food	55
##	10 usable	utensils	20	##	10 soils	clean	35	##	10 surfaces	food	55
##	11 utensils	improper	7	##	11 clean	exterior	34	##	11 labels	food	35
				##	12 clean	walls	32	##	12 food	utensil	25
				##	13 cooking	line	30	##	13 utensil	storage	25
				##	14 ceiling	tiles	29	##	14 food	ice	17
				##	15 provide	proper	28	##	15 food	thermometers	17

Performing statistical analysis on the chosen variables drew a better outline of the city and also helped me reinstate the value of these variables on my measure. The single variable t-test states that the license number seems to have a very feeble association with the violation level and Food related violation level, yet we see that it has positive incentive for violation level and yet holds

negative esteem with the rest of the factors. Food Related Violation is holding a similar negative connection with infringement level yet it has negative qualities with the license number too. It has 725 observational values though the 1000 perceptions. All the variables for the correlation coefficient of latitude and longitude is statistically significant, as these values are all less than 0.05 whereas the other variables are more towards having more chances of being observed at an error approximation.



All the issue related factors for the connection coefficient is measurably noteworthy, as these qualities are generally under 0.05 though the id variable is more towards having more odds of being seen at a blunder estimate. After analysing the above linear model we can say that the three \*\*\* mark at the PR value shows the best pairing with the Food related violation whereas the others seem to observe quite the contrary hence the only variable worth keeping for the model with respect to the violation level is food related violation.

We can observe from the Pearson's coefficient that food related violation has a strong correlation with the cleanliness related violation whereas the remaining two are not defined because of singularities and make less sense. We can observe from the Pearson's coefficient that cleanliness related violation has a strong correlation with the food related violation whereas the remaining



variable is not defined because of singularities and make less sense. We have perfect values of  $R^2$  and  $R$  hence stating that this model makes sense. If the beta coefficient is positive, the interpretation is that for every 1-unit increase in the predictor variable, the outcome variable will increase by the beta coefficient value and vice versa. In t-test, the null hypothesis is that the mean of the two samples is equal. This means that the alternative hypothesis for the test is that the difference of the mean is not equal to zero. In a hypothesis test, we want to reject or accept the null hypothesis with some confidence interval. Since we test the difference between the two means, the confidence interval in this case specifies the range of values within which the distinction may lie. The t-test will likewise deliver the p-value, which is the likelihood of wrongly dismissing the invalid speculation. The p-value is constantly contrasted and the essentialness dimension of the test. For cases, at 95% dimension of certainty, the critical dimension is 5% and the p-value is accounted for as  $p < 0.05$ . Small p-values propose that the invalid speculation is probably not going to be valid. The more smaller it is, we can confidently dismiss the invalid theory.

```

$r
      LICENSENO  ViolLevel FoodRelatedViol    lat.x
LICENSENO      1.00000000  0.04906644    -0.01828644 -0.26147821
ViolLevel      0.04906644  1.00000000    -0.16152513  0.08549021
FoodRelatedViol -0.01828644 -0.16152513  1.00000000  0.01751942
lat.x          -0.26147821  0.08549021  0.01751942  1.00000000
lon            -0.25947556  0.07726614  0.03484869  0.90893924
lon
LICENSENO      -0.25947556
ViolLevel      0.07726614
FoodRelatedViol 0.03484869
lat.x          0.90893924
lon            1.00000000

      LICENSENO      NA 0.1209950401117816376 0.5635389801025
ViolLevel      0.1209950401117816376 NA 0.0000002817828
FoodRelatedViol 0.5635389801024586198 0.0000002817828 NA
lat.x          0.0000000000008477663 0.0213273740968 0.6376785574913
lon            0.0000000000012803092 0.0375278306092 0.3487610652697
lon
lat.x
LICENSENO      0.0000000000008477663 0.000000000001280309
ViolLevel      0.0213273740967501180 0.037527830609222823
FoodRelatedViol 0.6376785574913090837 0.348761065269701742
lat.x          NA 0.000000000000000000
lon            0.000000000000000000 NA

```

Welch's two -t test states through the first set of values that the F value is 0.038, and p-value is high too. In other words, the variation of Violation level means with license status (numerator) is much smaller than the variation of Violation level within each status, and our p-value is greater

than 0.05 (as suggested by normal scientific standard). Hence we can conclude that for our confidence interval we reject the alternative hypothesis H1 that there is a significant relationship between LIC STATUS and Viol Level. The same was proven above in the t-test as well hence as we can see it connects and proves that the null hypothesis. In layman terms, we can see that it is true the Violation level is not relying on the status of the business chain which I have pointed out earlier.

The second set of F value is 4.09, and p-value is quite low too. In other words, the variation of Violation level with the Result of the health dept (numerator) is much greater than the variation of Violation level within each status, and our p-value is less than 0.05 (as suggested by normal scientific standard). Hence we can conclude that for our confidence interval we accept the alternative hypothesis H1 that there is a significant relationship between Result and Violation level. The Anova test accompanied with the Post hoc test was carried on to dive deeper into the dataset and point out distinctions.

```
Welch Two Sample t-test

data: ViolLevel by LICSTATUS
t = -0.19498, df = 12594, p-value = 0.8454
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -0.02570184  0.02105119
sample estimates:
mean in group 1 mean in group 3
 3.402156      3.404482
```

```
tuk<- TukeyHSD(anova4)
tuk

## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = ViolLevel ~ DESCRIPT, data = tracts_cleaned)
##
## $DESCRIPT
##
##           diff            lwr            upr
## Eating & Drinking w/ Take Out-Eating & Drinking -0.02821124 -0.062369658
## Mobile Food Walk On-Eating & Drinking 0.50868821 -0.079929059
## Retail Food-Eating & Drinking 0.04381387 0.007376908
## Mobile Food Walk On-Eating & Drinking w/ Take Out 0.53689944 -0.051779203
## Retail Food-Eating & Drinking w/ Take Out 0.07202510 0.034609676
## Retail Food-Mobile Food Walk On -0.46487434 -1.053689589
##
##           upr            p adj
## Eating & Drinking w/ Take Out-Eating & Drinking 0.005947187 0.1461161
## Mobile Food Walk On-Eating & Drinking 1.097305469 0.1177227
## Retail Food-Eating & Drinking 0.080250826 0.0108175
## Mobile Food Walk On-Eating & Drinking w/ Take Out 1.125578084 0.0884723
```

To figure out which groups are not quite the same as the others I directed a post hoc pair examination (note we can't play out numerous anova tests one for each pair, as this would build

our error) which is intended to assess pair means. There are many post hoc tests accessible for investigation of change and for my situation I utilised the “Tukey” post hoc test.

From the table above (taking a gander at diff and p-adj segments) I can see which description have noteworthy contrasts in Violation and food related infringement separately from others. For instance I inferred that, there is no huge distinction in food related infringement cases between any of the portrayal as all p esteems are more noteworthy than 0.05 while for infringement level it is enrolling no huge contrast Eating and Drinking w/Take Out-Eating and Drinking, Mobile Food Walk On-Eating and Drinking, Mobile Food Walk On-Eating and Drinking w/Take Out and Retail Food-Mobile Food Walk On. Although, there is a significant difference in violation level new cases between Retail Food-Eating & Drinking and Retail Food-Eating & Drinking w/ Take Out.



### Discussion:

In this paper, I have proposed a systematic procedure for accessing the violation level of various food establishments by constructing meaningful, science-based metrics by ranking results given out by the Health Department around the city. By using nonlinear urban scaling laws as a baseline, our procedure accounts for the underlying principles and socioeconomic dynamics that give rise to

cities to distinguish general effects of urbanism from local dynamics and, consequently, leads to a much simpler and direct perspective into the local factors that make or break specific places. We are hence able to layout our neighbourhoods in a distinguished way, the whole point of this paper is to use violation level as a parameter and enumerate the areas in the city accordingly.

From this perspective, the general measurably stable properties of urban areas develop as a pecking order of interrelated crucial amounts. To start with, Allston houses the least disregarded eateries contrasted with every single other territory of Boston. Hyde park envelops the most extreme eateries with fluctuations in ViolLevel. East Boston and Boston are viewed as two separate territories in the Boston city all through the dataset and East Boston has the most elevated middle of violLevel. Dorchester and Mattapan have dynamic authorised eateries. Unusually, the dynamic licenses have a high infringement status. FS, permit classification has greatest fizzles from the health department. Different restaurants under the same business names or food chains have different record for violation by health department.

Here, I have made the examination a stride further and demonstrated that the deviations of our selected variables from these conventional laws, which express neighbourhood factors explicit to singular urban communities, additionally show dispersions and relationships that are shockingly steady over long occasions as most of the inactive licenses still have the most violations and vice versa. It is accordingly uncommon that, regardless of the enormous decent variety of human and social conduct, the elements and association of urban frameworks, just as of individual urban areas, is an eminent unsurprising wonder. All out number of non-dynamic licenses are 3289, which implies that 47 % of the complete licenses offered out to the business chains are idle. All

out number of dynamic licenses are 3605, which implies that 52 % of the all out licenses offered out to the business chains are dynamic. The postal districts 2124, 2108, 2210 and 2128 are zones Ashmont, Beacon slope, central boston and Cambridge region separately. They appear to have every single dynamic permit eatery. This discloses to us that these areas are great spots to live in and the merchants take legitimate consideration of their licenses which demonstrates genuineness towards their calling, consequently this features these regions. The stick codes like 2116, 2127, 2110 and 2129 speak otherwise. The zones are Quincy, harbour end, north boston and Charlestown. Dorchester somehow maps under all of these issues, stating that it ain't a quaint neighbourhood. I observed that the food establishments that have cleanliness related issue definitely have animal and pest related issues as well we can see that through the graphs. We can state from the map above that the central boston has the most violated food chains compared to the other areas.

In synopsis, I have utilised the experimental signs of the basic standards of aggregation and the verifiable system structures for elements in charge of the development of urban communities to account efficiently for urban elements at various scales. This worldview enables us to isolate proportions of genuine nearby elements and association in urban communities from their nonexclusive widespread conduct. The whole measure which was laid forth to characterise the neighbourhood through the violation levels, violation descriptions and the health department's comments was achieved. I was able to envision a vivid picture of Boston where one can now locate the distinctions between different areas, the different description of food according to their license category, the most violated area etc.

Policies that focus instead on establishing beneficial fundamental change in local urban dynamics will be very difficult to achieve but very much worth creating, as they will position a city for a long run of prosperity and innovation. It would be interesting to investigate whether similar long term memory and persistence of urban dynamics is also a property of fast changing urban systems or not. Being able to illustrate and quantify these parameters will be essential in truly understanding how cities function, and ultimately will contribute to urban planning and developing policy solutions that will improve efficiency across the entire network, rather than to one area in isolation; we need to be able to assess the impact of changes on various scales, spatial and temporal, to create effective solutions in communities and for cities as whole entities. With the increasing availability of real-time and wide-scale urban data, I believe we should intend to explore some of these important questions in future work. From my perspective, this has the potential to really unpack the dynamics within and between different city systems and understand the effectiveness and evolution of these existing relationships when creating integrated and coordinated policies.