NYC Yellow Cab Demand Analysis

= 9G73

1091

Stephen Cho,

Minjee Kim

Introduction

- Objective: analyzing taxi demand and taxi traffic flow in New York City
- Dataset: Yellow taxi trip record data of Aug 2024 (provided by the <u>NYC TLC</u>)
- Data published on the TLC website, separated by year, month and vehicle type

NYC Taxi & Limousine Comm	mission				311 Search all NYC.go	v websites			Expand All	Collapse
		Taxi & Limousir	ne Commission	-	বাংল্য 🕨 Translate 🛛 🛡	Text-Size	▼ <u>2024</u>			
About Pas	ssengers Drivers	Vehicles	Businesses	TLC Online	Search	Q	January		July	
About TLC	Data and F	Reports	TLC Initi	atives	Contact TLC	<mark>;</mark>	 Yellow Taxi Tri Green Taxi Tri For-Hire Vehic (PARQUET) High Volume I Records (PAR 	rip Records (PARQUET) ip Records (PARQUET) cle Trip Records For-Hire Vehicle Trip RQUET)	 Yellow Taxi Trip Rec Green Taxi Trip Reco For-Hire Vehicle Trip (PARQUET) High Volume For-Hir Records (PARQUET) 	ords (PARQUE ords (PARQUE) Records re Vehicle Trip)
Pilot Programs	TLC Tri	p Reco	ord Data	a			February		August	
Reports	Yellow and green up and drop-off le reported passeng to the NYC Taxi a Taxicab & Livery I	Yellow and green taxi trip records include fields capturing pick-up and drop-off dates/times, pick- up and drop-off locations, trip distances, itemized fares, rate types, payment types, and driver- reported passenger counts. The data used in the attached datasets were collected and provided to the NYC Taxi and Limousine Commission (TLC) by technology providers authorized under the Taxicab & Livery Passenger Enhancement Programs (TPEP/LPEP). The trip data was not created					 Green Taxi Tri Green Taxi Tri For-Hire Vehic (PARQUET) High Volume I Records (PAR 	For-Hire Vehicle Trip	 Green Taxi Trip Reco Goren Taxi Trip Reco For-Hire Vehicle Trip (PARQUET) High Volume For-Hir Records (PARQUET) 	re Vehicle Trip
Reports	Yellow and green up and drop-off la reported passeng to the NYC Taxi a Taxicab & Livery I by the TLC, and T	taxi trip record ocations, trip di ger counts. The und Limousine (Passenger Enha TLC makes no r	s include fields ca istances, itemized data used in the a Commission (TLC) ancement Program representations as	pturing pick-up a fares, rate types ttached datasets by technology p ns (TPEP/LPEP). to the accuracy	and drop-off dates/time s, payment types, and c s were collected and p providers authorized un . The trip data was not of these data.	es, pick- lriver- rovided der the created	 Yellow Taxi Tri Green Taxi Tri For-Hire Vehic (PARQUET) High Volume I Records (PAR 	TIP Records (PARQUET) ip Records (PARQUET) cle Trip Records For-Hire Vehicle Trip RQUET)	 Yellow laxi Irip Rec Green Taxi Trip Recc For-Hire Vehicle Trip (PARQUET) High Volume For-Hir Records (PARQUET) 	or pro F re

Vehicle Types



- "Traditional" taxi (respond to street hails)
- More reliable data collection system (collected by TLC-authorized technology providers)



- Vehicles do not respond to street hails
- Data collected & provided by third-party corporations

Target narrowed down to Yellow and Green Taxi trip records

Vehicle Types (continued)



- Green taxi has very small trip counts (2% of total taxi trips)
- Mainly covers outer boroughs (cannot pick up new passengers in "yellow zone")
- Green Taxi trip record data does not fit our purpose, and is neglectable in size

\ /	-	- •	- •	
Yell	$\bigcap \setminus \Lambda /$	lavi	Irin	I)ata
		IUNI	μη	Data

	vendor_name	Trip_Pickup	DateTime	Trip_Dropoff_DateTime	Passenger_Coun	t Trip	_Distance <dbl></dbl>	Start_Lon <dbl></dbl>	Start_Lat <dbl></dbl>	Rate_Code +
1	VTS	2009-01-04	02:52:00	2009-01-04 03:02:00		1	2.63	-73.99196	40.72157	NA
2	VTS	2009-01-04	03:31:00	2009-01-04 03:38:00		3	4.55	-73.98210	40.73629	NA
3	VTS	2009-01-03	15:43:00	2009-01-03 15:57:00		5	10.35	-74.00259	40.73975	NA
4	DDS	2009-01-01	20:52:58	2009-01-01 21:14:00		1	5.00	-73.97427	40.79095	NA
5	DDS	2009-01-24	16:18:23	2009-01-24 16:24:56		1	0.40	-74.00158	40.71938	NA
						2	1 20	72 00001	40 72501	ALA.
6 6 r	DDS ows 1-9 of 18 c	2009-01-16 olumns	22:35:59	2009-01-16 22:43:35		2	1.20	-75.36361	40.75301	0 * ×
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward	22:35:59 End_Lon	2009-01-16 22:43:35 End_Lat Payment_Type	Fare_Amt s	urcharge	mta_tax	Tip_Amt	Tolls_Amt	Total_Amt
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbl></dbl>	22:35:59 End_Lon <db></db>	End_Lat <pre>cdbis</pre> Chrs	Fare_Amt s	urcharge	mta_tax <dbl></dbl>	Tip_Amt	Tolls_Amt	Total_Amt
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbl> NA</dbl>	22:35:59 End_Lon <dbl> -73.99380 72.95585</dbl>	2009-01-16 22:43:35 End_Lat Color= Color= Color= Color= End_Cash	Fare_Amt s	urcharge	mta_tax <dbl> NA</dbl>	Tip_Amt <dbl></dbl>	Tolls_Amt <dbl></dbl>	Total_Amt <dbl> 9.40</dbl>
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbl> NA NA</dbl>	22:35:59 End_Lon <dbl> -73.99380 -73.95585</dbl>	End_Lat Payment_Type <dbl> cchro 40.69592 CASH 40.76803 credit</dbl>	Fare_Amt s	urcharge <dbi> 0.5 0.5</dbi>	mta_tax <dbl> NA NA</dbl>	Tip_Amt <dbl> 0.00 2.00 4.74</dbl>	Tolls_Amt <dbl> 0</dbl>	Total_Amt <dbl> 9.40 14.60</dbl>
6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbl> NA NA NA</dbl>	22:35:59 End_Lon <dbl> .73.99380 .73.95585 .73.86998</dbl>	End_Lat Payment_Type <dbl> cchr> 40.69592 CASH 40.76803 Credit 40.77023 Credit</dbl>	Fare_Amt s	urcharge <dbi> 0.5 0.5 0.0 0.5 0.5 0.0 0.5 0.0 0.5 0.5</dbi>	mta_tax <dbl> NA NA</dbl>	Tip_Amt <dbl> 0.00 2.00 4.74</dbl>	Tolls_Amt <db>></db>	Total_Amt <dbl> 9.40 14.60 28.44</dbl>
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbi> NA NA NA NA NA</dbi>	22:35:59 End_Lon <dbl> -73.99380 -73.95585 -73.86998 -73.99656</dbl>	End_Lat Payment_Type <dbi> <chr> 40.69592 CASH 40.76803 Credit 40.77023 Credit 40.73185 CREDIT</chr></dbi>	Fare_Amt s <dbl> 8.9 12.1 23.7 14.9</dbl>	urcharge <dbl> 0.5 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.0 0.5 0.5</dbl>	mta_tax <dbl> NA NA NA NA</dbl>	Tip_Amt <dbl> 0.00 2.00 4.74 3.05</dbl>	Tolls_Amt <db > 0 0 0</db >	Total_Amt <dbi>9.40 14.60 28.44 18.45</dbi>
6 6 r	DDS ows 1-9 of 18 c store_an	2009-01-16 olumns d_forward <dbi> NA NA NA NA NA NA</dbi>	22:35:59 End_Lon <dbl> -73.99380 -73.95585 -73.86998 -73.99656 -74.00838</dbl>	End_Lat Payment_Type <dbl> cchr> 40.69592 CASH 40.77023 Credit 40.77023 Credit 40.73185 CREDIT 40.72035 CASH</dbl>	Fare_Amt s 8.9 12.1 23.7 14.9 3.7 14.9	urcharge <dbl> 0.5 0.5 0.0 0.5 0.5</dbl>	mta_tax <dbl> NA NA NA NA NA</dbl>	Tip_Amt <dbl> 0,00 2.00 4.74 3.05 0.00</dbl>	Tolls_Amt _ <dbl> 0 0 0 0</dbl>	Total_Amt <dbi>9.40 14.60 28.44 18.45 3.70</dbi>

Data Dictionary – Yellow Taxi Trip R	ecords May 11, 2022 Page 1 of 2					
This data distingant describes valle	u tavi trip data . For a dictionary describing green tavi data, or a man					
f the TLC Taxi Zones, please visit http://www.pyc.gov/html/tlc/html/about/trip, record, data.shtml						
Fire rec ran zones, prease visit <u>intp://www.nyc.gov/ittill/dout/thp_record_uata.sittill</u> .						
Field Name	Description					
VendorID	A code indicating the TPEP provider that provided the record.					
	1= Creative Mobile Technologies, LLC; 2= VeriFone Inc.					
tpep_pickup_datetime	The date and time when the meter was engaged.					
tpep_dropoff_datetime	The date and time when the meter was disengaged.					
Passenger_count	The number of passengers in the vehicle.					
	This is a driver-entered value.					
Trip_distance	The elapsed trip distance in miles reported by the taximeter.					
PULocationID	TLC Taxi Zone in which the taximeter was engaged					
DOLocationID	TLC Taxi Zone in which the taximeter was disengaged					
RateCodeID	The final rate code in effect at the end of the trip.					
	1= Standard rate					
	2=JFK					
	3=Newark					
	4=Nassau or Westchester					
	5=Negotiated fare					
	6=Group ride					
Store_and_fwd_flag	This flag indicates whether the trip record was held in vehicle					
	memory before sending to the vendor, aka "store and forward,"					
	because the vehicle did not have a connection to the server.					
	V- store and forward trip					
	I - store and forward trip					
	N= not a store and forward trip					

- Raw dataset (old version; left) consists of 18 columns, including key variables representing temporal (Pickup/Dropoff Date & Time) and spatial (Pickup/Dropoff coordinates) information.
- Each row is a yellow taxi trip record
- The TLC has replaced pickup/dropoff location details with "taxi zone" ID information for records since 2011
- Our goal is to analyze recent taxi demand patterns; need to work with the new format by generating (approximate) coordinates to perform spatial analysis

Yellow Taxi Data (continued)



1	OBJECTID	Shape_Leng	Shape_Area	zone	LocationID	borough	geometry	-
1	1	0.11635745	7.823068e-04	Newark Airport	1	EWR	list(list(c(933100.91835271, 93309	1.011480056, 933 [] 🔍
2	2	0.43346967	4.866340e-03	Jamaica Bay	1	Queens	list(list(c(1033269.24359129, 1033	439.64263915, 10 [] 🔍
3	3	0.08434111	3.144142e-04	Allerton/Pelham Gardens	1	Bronx	llst(llst(c(1026308.76950666, 1026	495.5934945, 102 [] 🔍
4	4	0.04356653	1.118719e-04	Alphabet City	4	4 Manhattan	list(list(c(992073.46679686, 99206	8.666992202, 992 [] ^Q .
5	5	0.09214649	4.979575e-04	Arden Heights	1	5 Staten Island	list(list(c(935843.310493261, 9360	46.564807966, 93 [] ^Q .
6	6	0.150/0054	6.0646104.04	Amorbar/East Wadoworth		States Island	Netvilletic/DEESER 74EEES761 DEEE	15 255504474 06 / 1 Q.
Show	ing 1 to 6 of 263	entries, 7 total colu	umns	Another in the should		Statemistand	asiasi(c)200000740002707,2000	1322334414, 30 [H]
Show	ing 1 to 6 of 263	entries, 7 total colu DOL	ocationID	PU_Longitude	PU_La	titude	DO_Longitude	DO_Latitude
Show	ling 1 to 6 of 263	entries, 7 total colu	ocationID	PU_Longitude	PU_La	titude <dbl></dbl>	DO_Longitude <dbl></dbl>	DO_Latitude <dbl></dbl>
Show	ling 1 to 6 of 263	entries, 7 total colu	ocationID <int> 161</int>	PU_Longitude <dbl> -73.96563 -73.98879</dbl>	PU_La 40.	titude <dbl> 76862 75351</dbl>	DO_Longitude <db> -73.97770 -73.99244</db>	DO_Latitude <dbl> 40.75803 40.74850</dbl>
PU	LocationID <int> 237 100 161</int>	entries, 7 total colu	ocationID <int> 161 186 114</int>	PU_Longitude <dbi> -73.96563 -73.98879 -73.9770</dbi>	PU_La 40. 40.	titude <dbi> 76862 75351 75803</dbi>	DO_Longitude <dbl> -73.97770 -73.99244 -73.99738</dbl>	DO_Latitude <dbl> 40.75803 40.74850 40.72834</dbl>
PU	locationID <int> 237 100 161 100</int>	entries, 7 total coli DOL	.ocationID <int> 161 186 114 13</int>	PU_Longitude <dbl> -73.96563 -73.98879 -73.97770 -73.98879</dbl>	PU_La 40. 40. 40. 40. 40.	titude <dbl> 76862 75351 75803 75351</dbl>	DO_Longitude <dbl> -73.97770 -73.99244 -73.99738 -74.01608</dbl>	DO_Latitude <dbl> 40.75803 40.74850 40.72834 40.71204</dbl>
PU	4 mg 1 to 6 of 263	entries, 7 total coli DOL	.ocationID <int> 161 186 114 13 75</int>	PU_Longitude <dbl> -73.96563 -73.98879 -73.97770 -73.98879 -73.98879 -73.94575</dbl>	PU_La 40. 40. 40. 40. 40. 40.	titude <dbl> 76862 75351 75803 75351 79001</dbl>	DO_Longitude <dbi> -73.97770 -73.99244 -73.99738 -74.01608 -73.94575</dbi>	DO_Latitude <dbl> 40.75803 40.74850 40.72834 40.71204 40.79001</dbl>

- The TLC also provides taxi zone details; great asset for calculating centroid coordinates and visualization
- NYC is divided into <u>263 taxi zones;</u> centroid coordinates are acceptable alternative for exact coordinates
- Columns have shape & geometric information, zone name, location ID, borough name
- PU_Longitude/Latitude, DO_Longitude/Latitude columns, each working as a pair, are mutated and merged to the Yellow Taxi Trip Dataset with PULocationID/DOLocationID used as reference

✓ All rows now have coordinates of pickup/dropoff locations

Data Cleaning



- Original dataset has 3 million rows; used <u>1 million</u> randomly extracted samples for efficiency
- Most variables often have missing/unusual values; only considered spatial & temporal variables for cleaning
- Spatial variables (PULocationID, DOLocationID) are intact for all rows; temporal variables more vulnerable
- "Duration": Gap between pickup and dropoff time (new variable); negative or extreme values removed
- "<u>trip_distance</u>": Extreme values removed (by IQR method)
- <u>859762</u> rows remain after data cleaning

Exploratory Data Analysis

- 1. Pickup/Dropoff Counts by Location
- 2. Pickup Counts by Time Periods

Pickup Counts by Taxi Zone



- Pickups are heavily focused in Manhattan borough, especially midtown Manhattan area
- 'Midtown Center' has the most pickups of 44892
- LaGuardia Airport and JFK Airport are the only two non-Manhattan area with significant volume of pickups
- Taxi zones in gray have no pickup recorded
- "Governor's Island/Ellis Island/Liberty Island" always have zero pickup counts since these areas can only be accessed by ferry boats

Dropoff Counts by Taxi Zone



- Dropoffs are also heavily focused in midtown Manhattan area, although slightly more spread out to other zones
- 'Midtown Center' also has the most dropoffs of 35758
- Outside Manhattan, dropoffs are less concentrated in LaGuardia Airport and JFK Airport
- Gray zones exist for dropoff counts as well

Pickup Counts by Hour



- 12 1PM has most pickup counts
- Pickup counts decline rapidly from 5 PM

Pickup Counts by Weekday



- Thursday has the most pickup counts
- Significantly less pickups on Mondays

Objective

- To understand the general traffic flow based on demands
- To capture the traffic flow from outer areas into the city during commuting hours and movements into the bar area during night times



Work Flow



Preparing the Data

8_Sun 🔅	7_Sun 🔅	6_Sun 🌐	5_Sun 🗧	4_Sun 🗦	3_Sun 🌐	2_Sun 🔅	1_Sun 🔅	0_Sun 🔅	DOLocationID
2	0	1	2	0	0	1	2	1	1
0	0	0	0	0	0	0	0	0	3
13	14	14	12	10	5	3	1	6	4
0	0	0	0	0	0	0	1	1	6
5	4	6	9	9	6	2	5	13	7
0	0	0	1	0	0	0	0	0	8
0	0	0	0	0	0	0	0	0	9
3	3	6	4	4	1	5	6	2	10
0	0	0	0	0	0	0	0	0	11
6	19	15	9	24	12	0	0	0	12
51	41	29	34	21	16	3	4	2	13
1	0	0	0	1	1	0	0	0	14
n	0	٥	٥	٥	n	0	1	1	15



Cyclic Time Behavior

Zone 88 in Manhattan



Cyclic Time Behavior

Zone 33 in Brooklyn



Fourier Transform on Time Series



(Lecture 7 – SpaceTime-Discrete)

A multi-resolution wavelet decomposition of a function $f_s(t)$ is an expression of the following form:

$$f_s(t) pprox eta_{s,00} \phi_{00}(t) + \sum_{j=-\infty}^J \sum_{k=0}^{2^{j-1}} eta_{s,j,k} \psi_{j,k}(t)$$

 $\beta_{s,00}$ is the scaling coefficient. The wavelets $\psi_{j,k}(t)$ are generated from a single wavelet $\psi(t)$, the so-called mother wavelet, by scaling and translation. The form of basis functions are known.

Fourier transform is a special case when $\psi(t) = e^{-2i\pi t}$

The temporal demand fs(t) can be represented by Fourier coefficients:

$$f_s(t) = eta_{s,0} + \sum_{k=1}^K eta_{s,k} \cos(2\pi k t) + \sum_{k=1}^K \gamma_{s,k} \sin(2\pi k t),$$

Our Model

$$Y(s,t)pprox f_s(t)+w_s+\epsilon(s,t)$$

Y(s,t): Observed demand (pickup counts) at location s and time t.

 $f_s(t)$: Temporal demand pattern at location s, capturing the temporal variability such as daily or weekly cycles.

 w_s : Spatial constraint, representing the inherent spatial connectivity.

 $\epsilon(s,t) :$ Error term capturing random noise or unmodeled variability.

$$f_s(t) = eta_{s,0} + \sum_{k=1}^K eta_{s,k} \cos(2\pi k t) + \sum_{k=1}^K \gamma_{s,k} \sin(2\pi k t) +$$

(Lecture 7 – SpaceTime-Discrete)

 $Y(s,t) = M(s,t)'\beta + w(s,t) + \epsilon(s,t)$

for $s \in D$ and $t \in [0, T]$. M(s, t) are local space-time covariate vectors β is an associated coefficient vector w(s, t): spatial temporal random effect. ϵ 's are pure error terms.

Use temporal basis $f_1(t), \dots, f_m(t)$: $w(s, t) = \sum_{i=1}^m f_i(t)\psi_i(s)$ we need to estimate spatially varying basis coefficients $\psi_i(s)$ (spatial functional data analysis)

Methodology (Chavent, 2017)

Aggregation measure based on the combined dissimilarity of two points i and j:



In matrix form:

 $\Delta_{\alpha} = (1 - \alpha)\Delta_0 + \alpha\Delta_1.$

- **D0:** captures how different two locations are based on their feature patterns
- D1: captures how geographically unconnected two locations are based on spatial adjacency
- When alpha = 0, the feature coefficients are clustered without any spatial smoothing

Standard Ward's Method: $I(\mathcal{C}_k) = \sum_{i \in \mathcal{C}_k} \sum_{j \in \mathcal{C}_k} \frac{w_i w_j}{2\mu_k} d_{ij}^2$ $I_{\alpha}(\mathcal{C}_k^{\alpha}) = (1 - \alpha) \sum_{i \in \mathcal{C}_k^{\alpha}} \sum_{j \in \mathcal{C}_k^{\alpha}} \frac{w_i w_j}{2\mu_k^{\alpha}} d_{0,ij}^2 + \alpha \sum_{i \in \mathcal{C}_k^{\alpha}} \sum_{j \in \mathcal{C}_k^{\alpha}} \frac{w_i w_j}{2\mu_k^{\alpha}} d_{1,ij}^2,$

Methodology (Chavent, 2017)

Model Component	ClustGeo Component	Description
$f_s(t)$: Temporal demand patterns	Feature-based dissimilarity $D_{ m 0}$	Temporal patterns (Fourier coefficients) are used to compute D_0 , capturing pairwise dissimilarities in temporal behavior across locations.
w_s : Spatial effect	Spatial dissimilarity D_1	Spatial relationships (from adjacency or proximity) are encoded in D_1 , penalizing clusters that split geographically connected locations.
$\epsilon(s,t)$: Random noise	Not explicitly modeled	ClustGeo assumes that noise is minor compared to the signal in D_{0} and D_{1} .
Combined effects	Combined dissimilarity Δ_lpha	$\Delta_lpha=(1-lpha)D_0+lpha D_1$ balances temporal and spatial effects in the clustering process.

Using ClustGeo in R

Hierarchical clustering with soft contiguity constraint.

The function hclustgeo implements a Ward-like hierarchical clustering algorithm with soft contiguity constraint. The main arguments of the function are:

- a matrix D0 with the dissimilarities in the "feature space" (here socio-economic variables for instance).
- a matrix D1 with the dissimilarities in the "constraint" space (here a matrix of geographical dissimilarities).
- a mixing parameter alpha between 0 an 1. The mixing parameter sets the importance of the constraint in the clustering procedure.
- a scaling parameter scale with a logical value. If TRUE the dissimilarity matrices D0 and D1 are scaled between 0 and 1 (that is divided by their maximum value).

The function choicealpha implements a procedure to help the user in the choice of a suitable value of the mixing parameter alpha.

Both hclustgeo and choicealpha can be combined to find a partition of the n = 303 French municipalities including geographical contiguity constraint. The two steps of the procedure are :

- 1. Find partition in K clusters of the 303 municipalities using the dissimilarity matrix D0. The clusters of this partition are homogeneous on the socio-economic variables and no contiguity constraint is used.
- 2. Choose a mixing parameter alpha in order to increases the geographical cohesion of the clusters (using the dissimilarity matrix D1) without deteriorating too much the homogeneity on the socio-economic variables.

https://cran.r-project.org/web/packages/ClustGeo/vignettes/intro_ClustGeo.html

ClustGeo in R

D0 "feature": time coefficients **D1 "constraint"**: geographical dissimilarities

Using ClustGeo:

- 1. Compute the Features
- 2. Compute the Spatial Constraints
- 3. Pick alpha
- 4. Cluster using hgeoclust()

Fourier Transform on Time Series

- Beneficial for capturing cyclical behaviors
- Steps
 - 1. Prepare the data
 - 2. Center the demand (helps the analysis focus on the deviation from the baseline)
 - 3. Use fft() in R to decompose the time series
 - 4. Extract the first four of real and imaginary components

Calculate the pairwise distance between the rows of scaled temporal data (used dist() function)

Heatmap of the scaled temporal coefficients 5 0 -5 00 Real1 Real2 Real3 Imag1 Imag4 lmag2 lmag3

Spatial Constraint

The spatial connectivity w_s is derived from k-nearest neighbor graph to define spatial connectivity.

 $w_{ij} = \begin{cases} 1 & \text{if locations } s_i \text{ and } s_j \text{ are spatially connected,} \\ 0 & \text{otherwise.} \end{cases}$

Calculate the spatial dissimilarity matrix as.dist(1-adj_matrix)

(Pick Up) Nearest Neighbor Graph, K = 5



Picking the best alpha

- Q0: Temporal homogeneity
- Q1: Spatial contiguity

Pick up Demand

Drop off Demand

K= 5 clusters

of 11%

based on D0

based on D1

1.0

0.8



alpha

0.6

0.4

Pick up Demand by Clusters





Pickup Cluster	Avg Trip Distance	Median Trip Distance	Avg Passenger Count	Avg Fare Amount	Avg Duration	Total Trips
1	0.271	0.00	1.791	71.233	0.882	201
2	3.237	3.27	1.373	22.088	14.061	36,348
3	1.915	1.51	1.294	13.358	12.007	375,302
4	1.817	1.49	1.368	13.480	12.838	348,398
5	2.127	1.79	1.410	14.157	12.952	99,513

Taxi Pickup Analysis: Time Effect





Drop off Demand by Clusters





Drop-off Cluster	Avg Trip Distance	Median Trip Distance	Avg Passenger Count	Avg Fare Amount	Avg Duration	Total Trips
1	0.387	0.000	1.761	68.022	1.692	216
2	3.656	3.850	1.341	22.481	16.476	52,203
3	1.761	1.470	1.374	13.325	12.631	414,877
4	2.453	2.000	1.331	15.069	12.932	142,211
5	1.641	1.370	1.292	12.293	11.346	250,255

Taxi Dropoff Analysis: Time Effect







Overall Traffic Flow: Matched Clusters





Traffic Flow: Into Manhattan





To analyze the traffic flow into Manhattan, we gathered the trips from <u>clusters 1, 2</u> (pick up) to <u>clusters 3,4,5</u> (drop off)

Traffic Flow: Into Manhattan





Traffic Flow: Out of Manhattan





To analyze the traffic flow out of Manhattan, we gathered the trips from <u>clusters 3,4,5</u> (pick up) to <u>clusters 1,2</u> (drop off)

Traffic Flow: Out of Manhattan





Traffic Flow: within Manhattan





For the traffic flow within Manhattan, we gathered the trips from (pick up) and to (drop off) <u>clusters 3,4,5</u>

Traffic Flow: within Manhattan





Traffic Flow Analysis

Pickup - Dropoff	Avg (SD) Trip Distance (miles)	Avg (SD) Passenger Count	Avg (SD) Fare Amount (\$)	Avg (SD) Duration (mins)	Total Trips
into city	4.47 (0.18)	1.32 (0.07)	26.62 (1.81)	21.09 (1.43)	5755
out of city	3.57 (1.83)	1.44 (0.51)	33.91 (14.74)	22.54 (2.17)	21625
Within city	2.16 (0.72)	1.36 (0.06)	14.41 (2.8)	13.17 (2.62)	801588

References

- New York City Taxi & Limousine Commission. (2024). *Yellow taxi trip data: August 2024*. Retrieved from https://www.nyc.gov/assets/tlc
- Strang, G. (2018). Fourier Transform. Massachusetts Institute of Technology. Retrieved November 19, 2024, from https://math.mit.edu/~gs/cse/websections/cse41.pdf
- Chavent, M., Kuentz-Simonet, V., Labenne, A., & Saracco, J. (2017). ClustGeo: An R package for hierarchical clustering with spatial constraints. *arXiv*. <u>https://doi.org/10.48550/arXiv.1707.03897</u>
- Chavent, M., Kuentz-Simonet, V., Labenne, A., & Saracco, J. (2021). Introduction to ClustGeo. The Comprehensive R Archive Network. Retrieved from https://cran.rproject.org/web/packages/ClustGeo/vignettes/intro_ClustGeo.html