

# **Transit demand estimation**

Pramesh Kumar

IIT Delhi

August 29, 2024

# Ridership

Usually measured using unlinked passenger trips or pass-km

# Significance

- ▶ give us an estimate of current and future transit needs
- ▶ important input for any service design
- ▶ help us select the best alternative among several alternatives at planning stage
- ▶ help us assess the effect of changes to the service, infrastructure, fares, etc.

# Factors affecting transit ridership

## Internal

- ▶ Fare
- ▶ Travel time (walking, waiting, transferring, in-vehicle)
- ▶ Service frequency
- ▶ Service coverage
- ▶ Stop location
- ▶ Route structure
- ▶ Transfers
- ▶ Comfort and convenience
- ▶ Information
- ▶ Crowding and reliability

## External

- ▶ Socio-economic factors (e.g., age, gender, income, auto ownership)
- ▶ Fuel prices
- ▶ Employment opportunities
- ▶ Land-use
- ▶ Safety
- ▶ Security
- ▶ Competition from other modes

**Remark.** Factors affecting passengers QoS also affect transit ridership. Better QoS help increase ridership.

## Forecasting techniques

- ▶ Expert judgment
- ▶ Rules of thumbs
- ▶ Surveys (e.g., stated preference)
- ▶ Elasticity
- ▶ Regression model
- ▶ Time series econometric model
  - Moving averages
  - Exponential smoothing
  - Double exponential smoothing (Holt's method)
- ▶ Trip distribution
- ▶ Discrete choice models
- ▶ Four step travel demand model

**Remark.** can vary based on scale (spatial, temporal) and market segments

## But ...

- ▶ Take forecasts with a pinch of salt
- ▶ Forecasts are usually wrong
- ▶ Aggregated forecasts are more accurate
- ▶ A good forecast is more than a single number
- ▶ Longer the forecast horizon, the less accurate the forecast will be

**Remark.** George Box said “All models are wrong but some are useful”

## Demand function and elasticity

Theoretically, demand  $D$  can be expressed as a function of various attributes (explanatory variables)  $x_1, \dots, x_m$ , i.e.,

$$D = f(x_1, \dots, x_m) \quad (1)$$

**Definition (Elasticity).** Percentage change in the demand wrt 1 % change in any attribute.

If  $D = f(x_i)$ , then

$$\epsilon_{D, x_i} = \frac{\frac{\Delta D}{D_0}}{\frac{\Delta x_i}{x_{i0}}} \quad (2)$$

**Example(s).** Simpson-Curtin rule: 3% fare increase reduces ridership by 1%

If  $\Delta x_i \rightarrow 0$ , then  $\epsilon_{D, x_i} = \frac{\partial D}{\partial x_i} \times \frac{x_{i0}}{D_0}$

# Elasticity

- ▶  $\epsilon_{D,x_i} < 0$  means demand curve is downward slopping, i.e., increase in  $x$  leads to decrease in the demand
- ▶  $\epsilon_{D,x_i} > 0$  means demand curve is upward slopping, i.e., increase in  $x$  leads to increase in the demand
- ▶  $\epsilon_{D,x_i} = 0$  means **perfectly inelastic demand**. This happens when there is no substitute for the current service.
- ▶  $|\epsilon_{D,x_i}| > 1$  means demand is elastic
- ▶  $|\epsilon_{D,x_i}| < 1$  means demand is inelastic
- ▶ Fare induces an inelastic demand.

**Remark.** In a competitive environment, a change in the attribute of one service may affect the demand of another service. Such changes are captured using **cross elasticity**.



# Regression modeling

1. State the problem
2. Model specification
  - An equation linking response and explanatory variables
  - Probability distribution of response variables
3. Parameter estimation
4. Check model adequacy
5. Inference

# Travel Demand Forecasting

We divide the geographical region into Transportation Analysis Zones (TAZs).

1. **Trip Generation** : Whether/when to travel? Estimates the number of trips from/to each zone.
2. **Trip Distribution** : Where to travel (which destination)? Estimates the other end of trips (OD trip matrix).
3. **Mode Choice** : How to travel (which mode)? Estimates the share of each mode from OD trips.
4. **Traffic Assignment** : How to travel (which route). Estimates traffic flow in transportation network.

## Trip distribution (OD estimation) methods

- ▶ No. of trips going from zone to another
- ▶ Expressed in the form of origin-destination passenger flow matrix
- ▶ Techniques
  - Growth factor method
  - Gravity method
  - Optimization
    - ▶ Entropy maximization
    - ▶ Maximum likelihood
    - ▶ Generalized least squares
  - Bayesian inference
  - Clustering
  - Trip chaining

# Growth factor method

Three types

1. Uniform
2. Singly-constrained
3. Doubly-constrained

Uniform

$$d_{\text{next year}}^{rs} = \gamma d_{\text{this year}}^{rs}, \forall (r, s) \in R \times S \quad (3)$$

where,  $\gamma$  is growth factor.

Issues

- ▶ Need to know the base demand which is not available for a new service
- ▶ All O-D pairs multiplied by the same growth factor. However, some areas can be developed more than others.

## Singly-constrained

- ▶ Origin-specific growth rate  $d_{\text{next year}}^{rs} = \gamma_r d_{\text{this year}}^{rs}, r \in R$
- ▶ Destination-specific growth rate  $d_{\text{next year}}^{rs} = \gamma_s d_{\text{this year}}^{rs}, s \in S$
- ▶ but not both

## Doubly-constrained

$$d_{\text{next year}}^{rs} = 0.5 \times (\gamma_r + \gamma_s) d_{\text{this year}}^{rs}, \forall (r, s) \in R \times S$$

If  $O_r \neq \sum_{s \in S} d^{rs}, \forall r \in R$  and  $D_s \neq \sum_{r \in R} d^{rs}, \forall s \in S$ , then we balance.

## Gravity model

- ▶ a widely-used, successful, aggregate model
- ▶ interaction between two locations:
  - increases with the amount of activity at each location
  - declines with increasing distance, time, and cost of travel between them
- ▶ general formula:

$$d^{rs} = \gamma_r \gamma_s O_r D_s f(c_{rs}) \quad (4)$$

- ▶ e.g., when the impedance is travel cost:

$$d^{rs} = \gamma_r \gamma_s \frac{O_r D_s}{c_{rs}} \quad (5)$$

If  $O_r \neq \sum_{s \in S} d^{rs}, \forall r \in R$  and  $D_s \neq \sum_{r \in R} d^{rs}, \forall s \in S$ , then we balance.

### Issues:

- ▶ Trip distribution and travel impedance are interdependent. Results of trip distribution should be used to update travel impedance.
- ▶ Does not take into account behavioral consideration. More sophisticated destination choice models that take into account user behavior in decision making should be used.

## Iterative proportional fitting (IPF)

1. Obtain the trips originated  $O_r$  (row sums) and destined  $D_s$  (column sums)
2. Obtain a seed matrix  $\{\hat{d}^{rs}\}_{(r,s) \in R \times S}$
3. Repeat the following steps:
  - $\hat{d}_{k+1}^{rs} = \frac{O_r}{\sum_{s \in S} \hat{d}_k^{rs}} \hat{d}_k^{rs}$ , where  $k$  is the iteration number.
  - $\hat{d}_{k+2}^{rs} = \frac{D_s}{\sum_{r \in R} \hat{d}_{k+1}^{rs}} \hat{d}_{k+1}^{rs}$
4. Repeat until  $\frac{O_r}{\sum_{s \in S} \hat{d}^{rs}}$  and  $\frac{D_s}{\sum_{r \in R} \hat{d}^{rs}} \approx 1$ .

### Issues:

- ▶ Non-structural zeros problem due to which a zero entry remains zero in every iteration.
- ▶ Quality of seed matrix should be good.

# Entropy maximization

## Notations

- ▶  $Z$ : set of zones
- ▶  $R$ : set of origins
- ▶  $S$ : set of destinations
- ▶  $d^{rs}$ : passenger trips from  $r$  to  $s$

From trip generation we know,

1. Trip generation

$$O_r = \sum_{s \in S} d^{rs}, \forall r \in R \quad (6)$$

2. Trip attraction

$$D_s = \sum_{r \in R} d^{rs}, \forall s \in S \quad (7)$$



We want to fill the following matrix

	$s_1$	$s_2$	.	.	.	$s_n$
$r_1$	$d^{r_1 s_1}$					
$r_2$		.				
.			.			
.				.		
$r_n$						$d^{r_n s_n}$

There are  $k = |R||S|$  entries in the OD matrix and total demand is  $Z = \sum_{r \in R} \sum_{s \in S} d^{rs}$ . Assuming that it is equally likely to travel on one of the  $k$  entries of the matrix. The probability that  $\{d^{rs}\}_{r \in R, s \in S}$  travelers will be traveling on individual O-D pairs is given by the multinomial probability distribution

$$\frac{Z!}{d^{r_1 s_1}! d^{r_2 s_2}! \dots d^k!} \left(\frac{1}{k}\right)^{d^{r_1 s_1}} \left(\frac{1}{k}\right)^{d^{r_2 s_2}} \dots \left(\frac{1}{k}\right)^{d^k}$$

$$= \frac{Z!}{d^{r_1 s_1}! d^{r_2 s_2}! \dots d^k!} \left(\frac{1}{k}\right)^Z$$

To maximize this, we take the logarithm

$$\begin{aligned}
 &= \log Z! - \sum_{(r,s) \in R \times S} \log d^{rs}! - Z \log k \\
 &= Z \log Z - Z - \sum_{(r,s) \in R \times S} (d^{rs} \log d^{rs} - d^{rs}) - Z \log k^1 \\
 &= \sum_{(r,s) \in R \times S} d^{rs} \log \left( \sum_{(r,s) \in R \times S} d^{rs} \right) - \sum_{(r,s) \in R \times S} d^{rs} \log d^{rs} - \left( \sum_{(r,s) \in R \times S} d^{rs} \right) \log k \\
 &= - \sum_{(r,s) \in R \times S} \frac{d^{rs}}{\sum_{(r,s) \in R \times S} d^{rs}} \left( \log \frac{d^{rs}}{\sum_{(r,s) \in R \times S} d^{rs}} \right) - \log k \\
 &= - \sum_{(r,s) \in R \times S} p^{rs} \log p^{rs} - \log k
 \end{aligned}$$

where,  $p^{rs}$  is the probability of traveling between  $(r, s) \in R \times S$

---

<sup>1</sup>Stirling's approximation  $\log x! \approx x \log x - x$

<sup>1</sup>Stirling's approximation  $\log x! \approx x \log x - x$

We get the following optimization problem

$$\begin{array}{ll} \underset{\mathbf{p}}{\text{maximize}} & - \sum_{(r,s) \in R \times S} (p^{rs} \log p^{rs}) \end{array} \quad (8a)$$

$$\begin{array}{ll} \text{subject to} & \sum_{s \in S} p^{rs} = \frac{O_r}{Z}, \forall r \in R \end{array} \quad (8b)$$

$$\sum_{r \in R} p^{rs} = \frac{D_s}{Z}, \forall s \in S \quad (8c)$$

$$p^{rs} \geq 0, \forall r \in R, \forall s \in S \quad (8d)$$

## Maximum likelihood estimation

- ▶ We assume that trips in OD pairs are i.i.d. random variables.
- ▶ Assuming Poisson distribution for the OD pairs, the probability of observing certain number of trips in that OD pair

$$\mathbb{P}(\hat{d}^{rs}) = \frac{(d^{rs})^{\hat{d}^{rs}}}{\hat{d}^{rs}!} e^{-d^{rs}}$$

where,  $d^{rs}$  is the estimated number of trips and  $\hat{d}^{rs}$  is the trips in the seed matrix.

- ▶ The likelihood function is given by

$$L = \prod_{(r,s) \in R \times S} \frac{(d^{rs})^{\hat{d}^{rs}}}{\hat{d}^{rs}!} e^{-d^{rs}}$$

We get the following optimization problem

$$\underset{\mathbf{d}}{\text{maximize}} \quad \log L \quad (9a)$$

$$\text{subject to} \quad \sum_{s \in S} d^{rs} = O_r, \forall r \in R \quad (9b)$$

$$\sum_{r \in R} d^{rs} = D_s, \forall s \in S \quad (9c)$$

$$d^{rs} \geq 0, \forall r \in R, \forall s \in S \quad (9d)$$

## Generalized least squares

Let us express the conservation constraints (using on-off counts from APC data or link counts from ETM data) as  $A\mathbf{d} = \mathbf{b}$ , where,

$$\mathbf{d} = \{d^{rs}\}_{(r,s) \in R \times S}.$$

$$\underset{\mathbf{d}}{\text{minimize}} \quad (A\mathbf{d} - \mathbf{b})^T W^{-1} (A\mathbf{d} - \mathbf{b}) + (\mathbf{d} - \hat{\mathbf{d}})^T V^{-1} (\mathbf{d} - \hat{\mathbf{d}}) \quad (10a)$$

- ▶  $W, V$  are weighting matrices (typically diagonal).
- ▶ The second term is referred to as a regulariser. Regularisation make sure that the estimated OD is not significantly deviating from the seed OD matrix.

## Bayesian estimation

- Bayes' theorem gives the posterior of unknown parameters (trips in the OD matrix)  $\theta$  given an observed measurement  $Y$  (link or on-off counts) as proportional to the likelihood of the observation and prior probability of the unknowns  $p(\theta)$

$$\mathbb{P}(\theta|Y) \propto \mathbb{P}(Y|\theta)\mathbb{P}(\theta)$$

- Estimates can be obtained as those giving the maximum a posteriori (MAP) density.

## Trip chaining

- ▶ AFC systems can be of two types:
  - Open: Only passengers' boarding/alighting location is recorded (usually transit systems with fixed fare)
  - Closed: Passengers' both boarding and alighting locations are recorded (usually transit systems with distance-based fare)
- ▶ In case of open system, alighting locations in the AFC data are inferred based on **rule-based** heuristics by making use of schedule, AVL and/or APC data.
  - Rule-based trajectories are usually based on walking time threshold, waiting time threshold, and space-time constraints.
- ▶ In both cases, transfers also need to be inferred in order to get entire trajectory.



# Clustering

- ▶ When boarding stop is not available in the AFC data, then clustering can be used to assign the boarding GPS locations to various transit stops.
- ▶ Clustering methods
  - K-means clustering
  - DBSCAN
  - others

Thank you!