
Aggregation and forecasting

Michel Bierlaire

`michel.bierlaire@epfl.ch`

Transport and Mobility Laboratory

Aggregation

- So far, prediction of individual behavior
- In practice, not useful
- Need for forecast of aggregate demand:
 - number of trips
 - number of passengers
 - etc.

Aggregation

Linear models

$$h_n = \alpha + \beta y_n$$

where

- h_n : quantity of energy n consumed
- y_n : price of energy n
- If \bar{y} is the average price
- $\bar{h} = \alpha + \beta \bar{y}$ is the average consumption

It does not work with choice models, because they are nonlinear

Aggregation

- “Travel/no travel” model, y_n income

$$\text{No travel } V_1 = 0$$

$$\text{Travel } V_2 = -3 + 3y_n$$

| | Income | V1 | V2 | P1 | P2 |
|--------------------|--------|----|------|-----|------|
| Household 1 | 1 | 0 | 0 | 50% | 50% |
| Household 2 | 10 | 0 | 27 | 0% | 100% |
| Avg. income | 5.5 | 0 | 13.5 | 0% | 100% |
| Avg. probabilities | | | | 25% | 75% |

Aggregation

- Choice model

$$P(i|x_n)$$

where x_n gathers attributes of all alternatives and socio-economic characteristics of n

- If the population is composed of N individuals, the total expected number of individuals choosing i is

$$N(i) = \sum_{n=1}^N P(i|x_n)$$

- Hopeless to know x_n for every and each individual
- The sum would involve a lot of terms.
- The distribution of x could be used.

Aggregation

- Assume that the distribution of x is continuous with PDF $p(x)$
- Then the share of the population choosing i is given by

$$\widehat{W}(i) = \int_x P(i|x)p(x)dx$$

- In practice, $p(x)$ is also unknown
- The integral may be cumbersome to compute

Aggregation

Most practical method: **sample enumeration**

- Population is assumed to be segmented into homogenous segments
- Let n be an observation in the sample belonging to segment s
- Let W_s be the weight of segment s , that is

$$W_s = \frac{N_s}{S_s} = \frac{\# \text{ persons in segment } s \text{ in population}}{\# \text{ persons in segment } s \text{ in sample}}$$

- The number of persons choosing alt. i is estimated by

$$\hat{N}(i) = \sum_{n \in \text{sample}} \sum_s W_s P(i|x_n) I_{ns}$$

where $I_{ns} = 1$ if individual n belongs to segment s , 0 otherwise

Aggregation

We can write

$$\begin{aligned}\hat{N}(i) &= \sum_{n \in \text{sample}} \sum_s W_s P(i|x_n) I_{ns} \\ &= \sum_{n \in \text{sample}} P(i|x_n) \sum_s W_s I_{ns}\end{aligned}$$

The term $\sum_s W_s I_{ns}$ is the weight W_n of individuals n belonging to segment s .

The **share** of alt. i is estimated by $W(i) =$

$$\frac{1}{N} \sum_{n \in \text{sample}} P(i|x_n) W_n = \frac{1}{N} \sum_{n \in \text{sample}} P(i|x_n) W_n$$

Illustration

The travel model:

- “Travel/no travel” model, y_n income

$$P(\text{travel}) = \frac{e^{-3+3y_n}}{1 + e^{-3+3y_n}}$$

- Population: $N = 200'000$ persons
- Sample: $S = 500$ persons
- Sampling rate: $S/N = 1/400$

Illustration

| s | y_s | S_s | N_s | $P(\text{travel})$ | PS_s | PN_s |
|-----|-------|-------|--------|--------------------|--------|--------|
| 1 | 0 | 150 | 20000 | 4.7% | 7 | 949 |
| 2 | 0.5 | 200 | 30000 | 18.2% | 36 | 5473 |
| 3 | 1 | 40 | 50000 | 50.0% | 20 | 25000 |
| 4 | 1.5 | 10 | 50000 | 81.8% | 8 | 40879 |
| 5 | 2 | 50 | 30000 | 95.3% | 48 | 28577 |
| 6 | 2.5 | 50 | 20000 | 98.9% | 49 | 19780 |
| | | 500 | 200000 | | 169 | 120657 |

$$120657 \neq 400 \times 169 = 67542$$

People with low probability of travel are oversampled

Forecasting

- Modify x_n in the sample to reflect anticipated modifications
- Apply the sample enumeration again
- Examples:
 - Socio-economic characteristics: scenarios of future demographics (level of education, modification of incomes, etc.)
 - Attributes of alternatives:
 - Policy variables: variables that we control (price, level of service, etc.)
 - Scenarios: scenarios about the competition

Example: characteristics

| s | y_s | S_s | P(travel) | W_s | Trips |
|-----|-------|-------|-----------|--------|--------|
| 1 | 0 | 150 | 4.74% | 133.33 | 949 |
| 2 | 0.5 | 200 | 18.24% | 150 | 5473 |
| 3 | 1 | 40 | 50.00% | 1250 | 25000 |
| 4 | 1.5 | 10 | 81.76% | 5000 | 40879 |
| 5 | 2 | 50 | 95.26% | 600 | 28577 |
| 6 | 2.5 | 50 | 98.90% | 400 | 19780 |
| | | | | | 120657 |

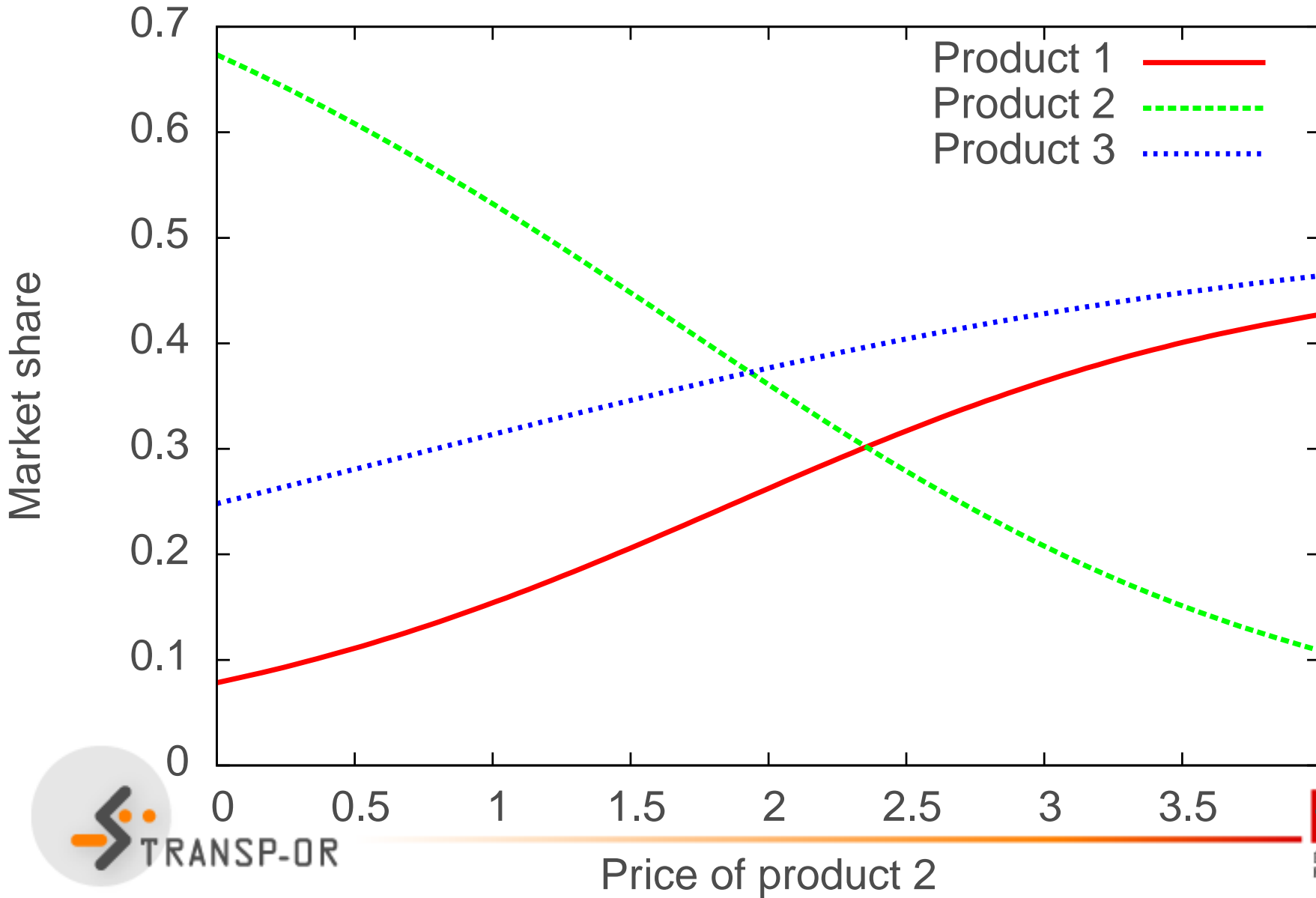
- Increase all salaries by 0.5
- What is the impact on the total number of trips?

Example: characteristics

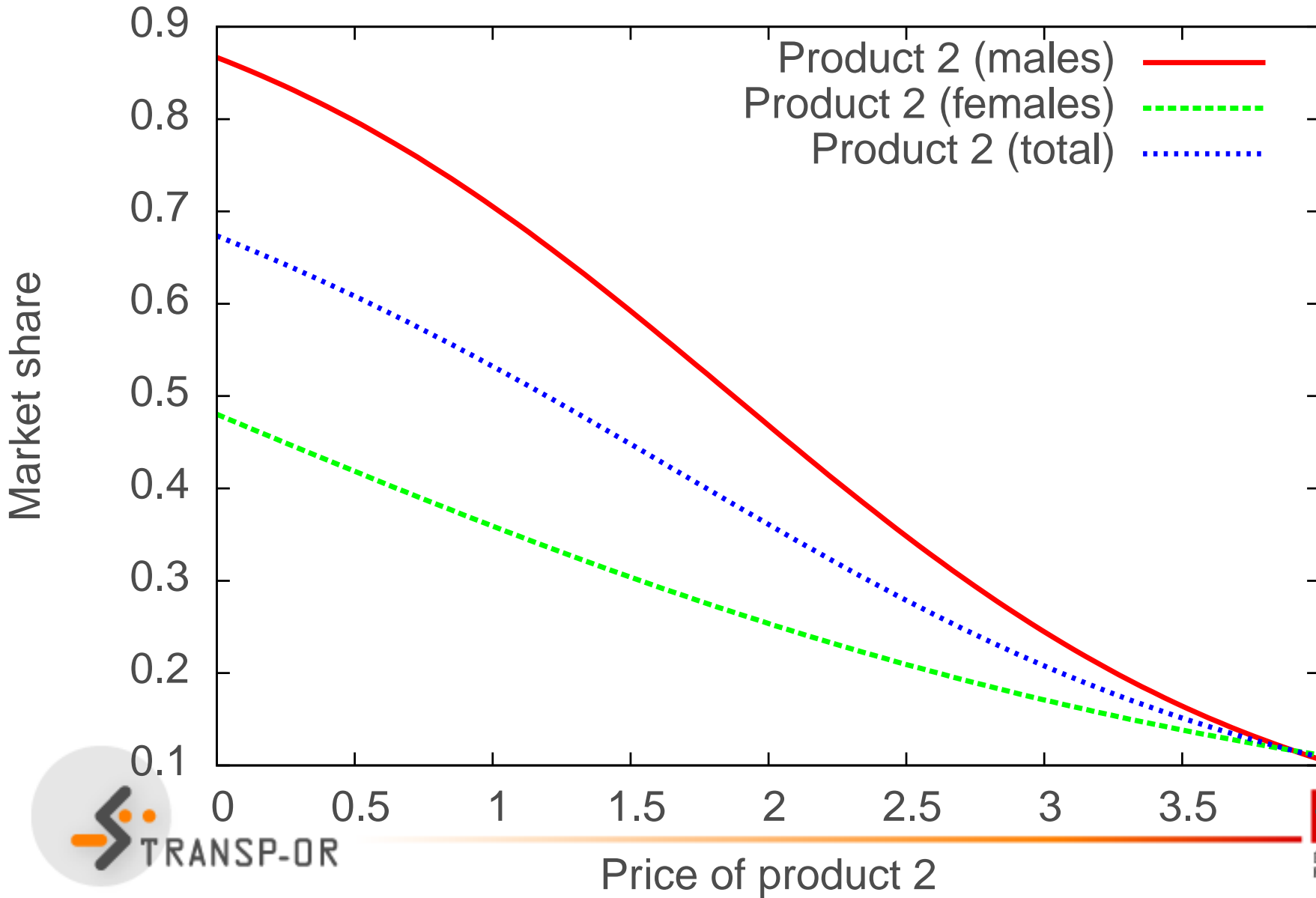
| s | y_s | S_s | P(travel) | W_s | Trips |
|-----|-------|-------|-----------|--------|--------|
| 1 | 0.5 | 150 | 18.24% | 133.33 | 3649 |
| 2 | 1 | 200 | 50.00% | 150 | 15000 |
| 3 | 1.5 | 40 | 81.76% | 1250 | 40879 |
| 4 | 2 | 10 | 95.26% | 5000 | 47629 |
| 5 | 2.5 | 50 | 98.90% | 600 | 29670 |
| 6 | 3 | 50 | 99.75% | 400 | 19951 |
| | | | | | 156777 |

- Before: 120657
- After: 156777
- Increase: about 30%

Example: attributes



Example: attributes



Price optimization

- Choice model captures demand
- Demand is elastic to price
- Predicted demand varies with price, if it is a variable of the model
- In principle, the probability to use/purchase an alternative decreases if the price increases.
- The revenue per user increases if the price increases.
- Question: what is the optimal price to optimize revenue?

In short:

- Price \uparrow \Rightarrow profit/customer \uparrow and number of customers \downarrow
- Price \downarrow \Rightarrow profit/customer \downarrow and number of customers \uparrow
- What is the best trade-off?

Revenue calculation

Number of persons choosing alternative i in the population

$$\hat{N}(i) = \sum_{n \in \text{sample}} P(i|x_n, p_{in})W_n$$

where

- p_{in} is the price of item i for individual n
- x_n gathers all other variables corresponding to individual n
- $P(i|x_n, p_{in})$ is the choice model
- W_n is the weight of individual n .

Revenue calculation

The total revenue from i is therefore:

$$R_i = \sum_{n \in \text{sample}} W_n P(i | x_n, p_{in}) p_{in}$$

If the price is constant across individuals, we have

$$R_i = p_i \sum_{n \in \text{sample}} W_n P(i | x_n, p_i)$$

Price optimization

Optimizing the price of product i is solving the problem

$$\max_{p_i} p_i \sum_{n \in \text{sample}} W_n P(i|x_n, p_i)$$

Notes:

- It assumes that everything else is equal
- In practice, it is likely that the competition will also adjust the prices

Illustrative example

A binary logit model with

$$V_1 = \beta_p p_1 - 0.5$$

$$V_2 = \beta_p p_2$$

so that

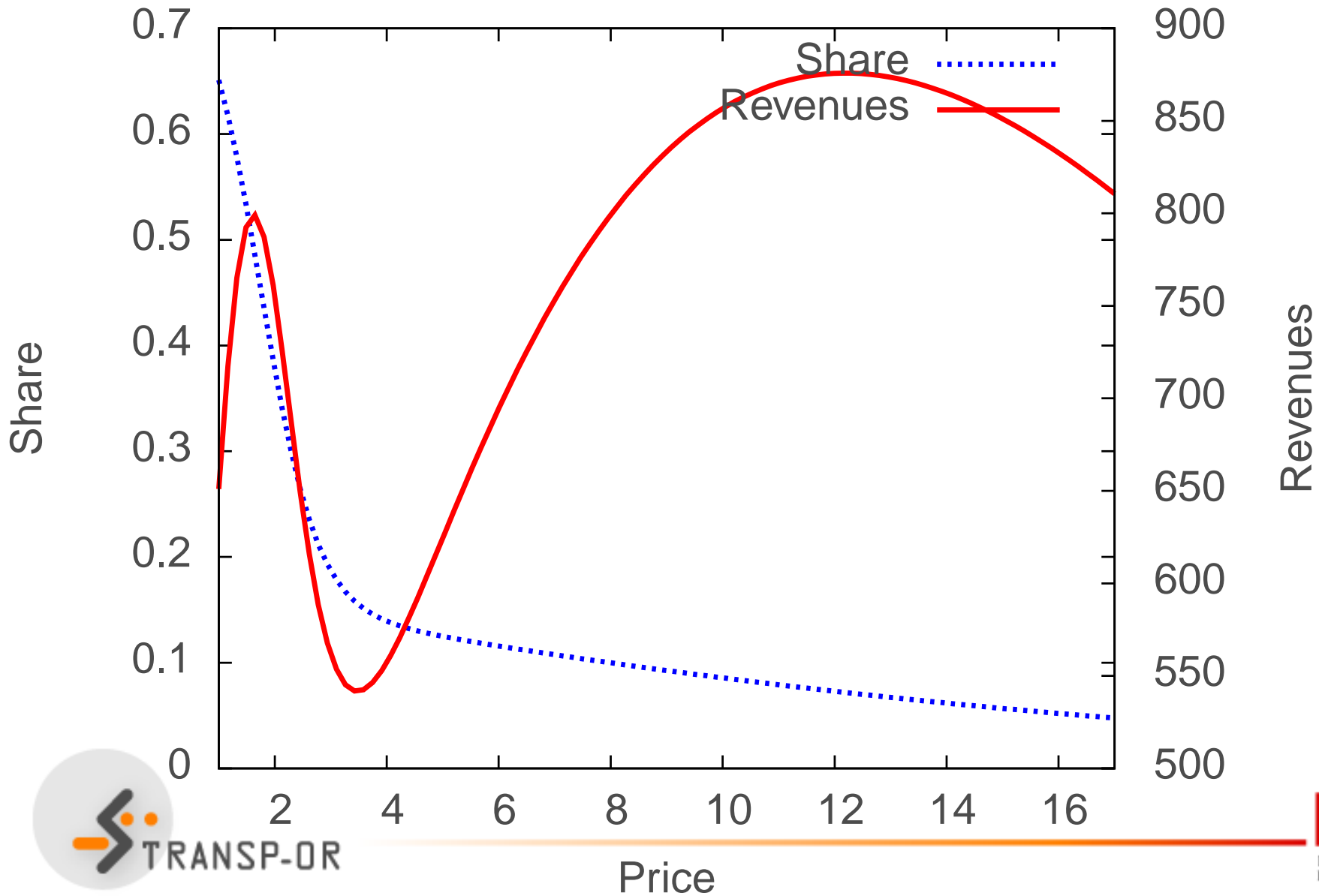
$$P(1|p) = \frac{e^{\beta_p p_1 - 0.5}}{e^{\beta_p p_1 - 0.5} + e^{\beta_p p_2}}$$

Two groups in the population:

- Group 1: $\beta_p = -2$, $N_s = 600$
- Group 2: $\beta_p = -0.1$, $N_s = 400$

Assume that $p_2 = 2$.

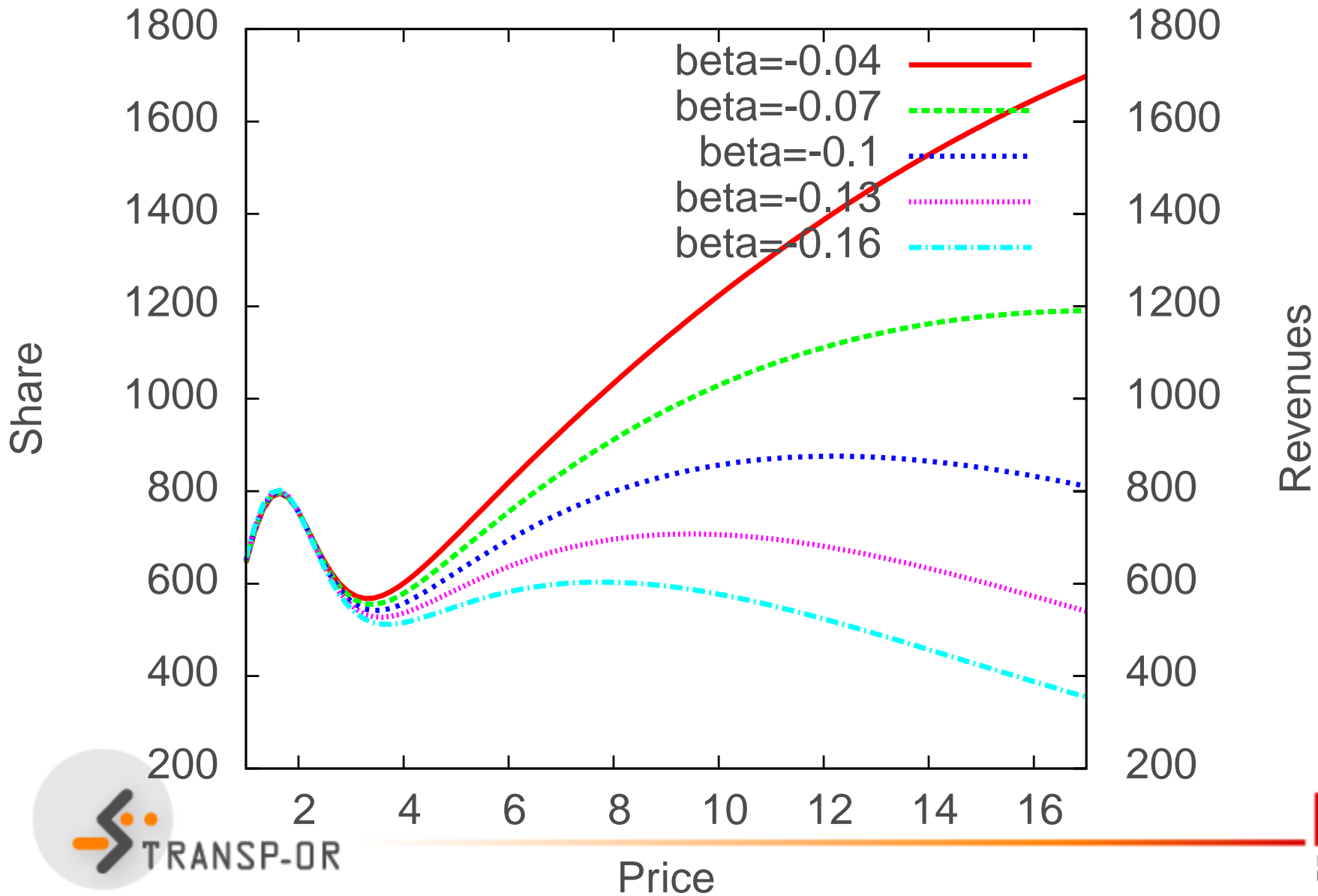
Illustrative example



Sensitivity analysis

- Parameters are estimated, we do not know the real value
- 95% confidence interval: $[\hat{\beta}_p - 1.96\sigma, \hat{\beta}_p + 1.96\sigma]$
- Perform a sensitivity analysis for β_p in group 2

Sensitivity analysis



Comments

- The estimation sample already contains a sample of the variables x .
- It is convenient to use the same sample for estimation and sample enumeration.
- It is valid only if revealed preference data (i.e. revealing real behavior) is used.
- Stated preference data (i.e. choice based on hypothetical situation) cannot be used for sample enumeration.