

Estimating the preference for safety features in cars

Tool: Multinomial logit & Mixed logit

The Analytics Edge

Companies such as General Motors carry out conjoint studies with customers to understand the tradeoff (valuation) of different attributes that make up a product or a service. Using data on preferences for safety features in new vehicles and then building models of discrete choice, the company can obtain estimates on the valuation of attributes. This provides important information that can be used to infer the effect of introducing new products in the market, pricing the product and designing the product features.

The models incorporate the comparison of attributes across alternatives and can also be used to capture heterogeneity in the choice modeling process.

Which of the following packages would you prefer the most?

Choose by clicking one of the buttons below

<p>Full Speed Range Adaptive Cruise Control</p> <p>-</p> <p>Traditional Back Up Aid</p> <p>Lane Departure Warning</p> <p>Front Collision Warning</p> <p>Side Head Air Bags</p> <p>Emergency Notification with Pictures</p> <p>-</p> <p>-</p> <p>Option Package Price: \$3,000</p> <p><input type="radio"/></p>	<p>Adaptive Cruise Control</p> <p>Navigation System with Curve Notification & Speed Advisor</p> <p>Rear Vision System</p> <p>-</p> <p>-</p> <p>Side Body Air Bags</p> <p>Emergency Notification</p> <p>-</p> <p>Head Up Display</p> <p>Option Package Price: \$500</p> <p><input type="radio"/></p>	<p>Full Speed Range Adaptive Cruise Control</p> <p>Traditional Navigation System</p> <p>-</p> <p>-</p> <p>Front Collision Warning</p> <p>Side Body & Head Air Bags</p> <p>-</p> <p>Night Vision with Pedestrian Detection</p> <p>Head Up Display</p> <p>Option Package Price: \$12,000</p> <p><input type="radio"/></p>	<p>NONE: I wouldn't purchase any of these packages</p> <p><input type="radio"/></p>
--	---	---	---

Analytics on safety feature data: R

Safety ← read.csv("safety data.csv")

Stor (safety)

9500 observations of

Summary (safety)

112 variables

head (safety)

Each customer performs

19 choice tasks with

a total of 500 customers

Case = Customer number (1 to 500)

No = Observation number (1 to 9500)

Task = Task number for a customer (1 to 19)

Each customer has a choice among four

packages in each choice task. (Stated preference data set)

For example customer 1 in the first choice task & chooses Ch2

Alternative 1
(Ch 1)

Alternative 2
Ch 2

Alternative 3
(Ch 3)

Alternative 4
(Ch 4)

GN 1

NS 4

BU 6

FP 3

RP 1

PP 1

TS 1

NV 1

MA 4

Price 2

level
of
attributes
in
package 1

GN 2

NS 1

BU 2

FP 2

RP 2

PP 1

TS 1

NV 1

MA 1

Price 2

Package
2

GN 2

NS 4

BU 5

FP 3

RP 3

PP 3

TS 3

NV 1

MA 1

Price 2

Package
3

All
at
zero
level

0

0

0

0

0

0

0

0

0

0

0

CC	Cruise control	3 levels
GN	Go notifier	2 levels
NS	Navigation system	5 levels
BU	Backup aids	6 levels
FA	Front park assist	2 levels
LD	Lane departure	3 levels
BZ	Blind zone alert	3 levels
FC	Front collision warning	2 levels
FP	Front collision protection	4 levels
RP	Rear collision protection	2 levels
PP	Parallel park aids	3 levels
KA	Knee air bags	2 levels
SC	Side airbags	4 levels
TS	Emergency notification	3 levels
NV	Night vision system	3 levels
MA	Driver assisted adjustment	4 levels
LB	Low speed braking assist	4 levels
AF	Adaptive front lighting	3 levels
HU	Head up display	2 levels
Price	Price	11 levels

(\$500, 1000, 1500, 2000, 2500,
3000, 4000, 5000, 7500,
10000, 12000)

The variables from "segment" to "income" give demographic information

Ch 1 = 1 if package 1 is chosen, 0 o/w

Ch 2 = 1 if package 2 is chosen, 0 o/w

Ch 3 = 1 if package 3 is chosen, 0 o/w

Ch 4 = 1 if no package is chosen, 0 o/w

Choice = $\begin{cases} \text{Ch 1} \\ \text{Ch 2} \\ \text{Ch 3} \\ \text{Ch 4} \end{cases}$ Depending on which option is chosen

table (Safety \$ Choice)	Ch 1	Ch 2	Ch 3	Ch 4
	2165	2404	2136	2795

Segment = Segment of car population that individual has

Year = Year

Miles = Mileage

Night = % of time customer drives at night

Gender = Gender

Gender = Gender (Female or Male)

Age = Age bracket

Educ = Education level (^{4 year} college, grade school, high school, postgrad, college, voc ^{1-3 years})

region = resident of region (NW, NE, SE, SW, W)

Urban = resident of Rural, Suburban, Urban

Income = Income

which (colnames(safety) == "CC1") 4

which (colnames(safety) == "Price 4") 83

Columns 4 to 83 of the safety dataframe contain the attribute levels for each of the 4 alternatives (choices) indexed by 1, 2, 3 and 4 respectively. Note the zero level in all cases indicate the attribute is not available with the fourth choice being the NONE option.

In general higher levels for attributes correspond to more technically advanced features.

We can capture these variables either directly or by adding dummy variables that take a value of 1 or 0 depending on the level.

Developing a logit model for the data

```
S <- mlogit.data (subset(safety, Task ≤ 12),  
  shape = "wide", choice = "Choice",  
  varying = c(4:83), sep = " ",  
  alt.levels = c("Ch1", "Ch2", "Ch3", "Ch4"),  
  id.var = "Case")
```

This command takes the dataframe where we use the first 12 choice tests for each customer in our training set (12 out of 19), indicates that the shape of the dataframe is wide (where each row is an observation), indicates that the variable indicating the choice made is defined by "Choice", the separator of the variable name and the alternative name (this helps to guess variables and alternative names), `varying = c(4:83)` (helps indicate the variables indices that are alternative specific), `alt.levels` indicates the names of the alternatives and `id.var` indicates the name of the variable containing the individual index (here "Case")

head(S)
str(S)

This creates a data frame of
the long format to which we
can use the mlogit package
(24000 = 6000 × 4 observations)

M ← mlogit (Choice ~ CC + GN + NS + BU + FA + LD +
BT + FC + FP + RP + PP + KA
+ SC + TS + NV + MA + LB + AF
+ MU + Price - 1, data = S)

Summary(M)

Log-likelihood $LL(\hat{\beta}) = -7567.9$
at optimality

$$AIC = -2LL(\hat{\beta}) + 2\left(\frac{\text{par}}{\text{number}}\right) = 15175.8$$

$$\begin{aligned}\text{Likelihood ratio} &= \frac{LL(\hat{\beta})}{LL(0)} \\ &= 1 - \left(\frac{-7567.9}{6000 \log\left(\frac{1}{4}\right)} \right) \\ &= 1 - \left(\frac{-7567.9}{-8317.76} \right) = 0.09\end{aligned}$$

Results indicate that

CC (cruise control), KA (knee air bags),

TS (emergency notification), MA (driver adjusted assessment), Price are very significant at 0.001 level.

As should be expected Price has a negative Coefficient.

LD (lane departure) and SC (side air bags) are significant at 0.01 level.

The non price coefficients have positive signs indicating that these are valued (higher levels of safety features).

To check the correct production in sample simply used on the maximum predicted choice probability

Actual Choice \leftarrow Subset (safety, Task ≤ 12)[, "Choice"]

P \leftarrow predict(M, newdata = S)

Predicted Choice \leftarrow apply(P, 1, which.max)

Actual choice keeps track of actual observed choice

P predicts the choice probability for the given dataset

The apply function applies to the matrix P,

across rows (indexed by 1), the function of

which.max (identifies index that is maximum)

Willingness to pay

Ratio of β_j coefficients can be used to estimate the willingness to pay for a particular attribute

Suppose β_1 is the coefficient for attribute 1 and β_2 is the coefficient of attribute 2 (price)

Note that β_2 will be typically negative.

Since more the price lesser the utility.

Assume say β_1 is positive.

$$U = \beta_1 x_1 + \beta_2 x_2 + \dots$$

Assume unit change (increase) in attribute 1

$$U = \beta_1 (x_1 + 1) + \beta_2 (x_2 + \Delta) + \dots$$

$$\text{Then } \Delta = - \frac{\beta_1}{\beta_2} \quad \left(\begin{array}{l} \text{Willingness to} \\ \text{pay for unit} \\ \text{increase in attribute} \end{array} \right)$$

Example, Willingness to pay for CC.

Cruise control.

$$\text{WTP} = \frac{0.1085}{0.2036} = 0.5329$$

Table (Predicted Choice, Actual Choice)

	Ch1	Ch2	Ch3	Ch4	
1	616	273	226	356	Predicted
2	222	647	235	320	
3	194	258	629	309	
4	339	352	324	700	
	Actual				

$$\text{Correct predictions} = \frac{616 + 647 + 629 + 700}{6000} = \boxed{0.432}$$

Assuming a complete random choice for predicting the choice would result in $\boxed{0.25}$.

```
Test <- mlogit.data(subset(safety, Test > 12),
  shape = "wide", choice = "Choice",
  varying = c(4:83), sep = "", alt.levels =
  c("Ch1", "Ch2", "Ch3", "Ch4"), id.var = "Case")
```

```
Test.Predict <- predict(M, newdata = Test)
```

```
Actual Choice <- subset(safety, Test > 12)[, "Choice"]
```

```
Predicted Choice <- apply(Test.Predict, 1, which.max)
```

	Ch1	Ch2	Ch3	Ch4
1	376	120	93	176
2	119	435	97	227
3	118	124	356	190
4	181	193	176	512

Table (Predicted Choice, Actual Choice)

$$\text{Correct predictions} = \frac{(376 + 435 + 356) + 512}{3500} = \boxed{0.4811}$$

Mixed logit model

To capture additional features that multinomial logit cannot capture such as random taste variation.

Standard logit
$$\tilde{U}_{ik} = \beta' x_{ik} + \tilde{\epsilon}_{ik}$$

In mixed logit, $\tilde{\beta}$ is modeled as a random parameter.

$$\tilde{U}_{ik} = \tilde{\beta}' x_{ik} + \tilde{\epsilon}_{ik}$$

Standard logit

$$P(Y_i = k) = \frac{e^{\beta' x_{ik}}}{\sum_{l=1}^K e^{\beta' x_{il}}}$$

Mixed logit

$$P(Y_i = k) = \int \frac{e^{\beta' x_{ik}} f(\beta) d\beta}{\sum_{l=1}^K e^{\beta' x_{il}}}$$

where $f(\beta)$ is the density function of $\tilde{\beta}$.

Integral of logit probabilities

Mixed logit is computationally more challenging to solve due to the use of simulation.

Optimization methods.

The problems are no longer convex in this setting and finding a global optimum might not

From an estimation perspective, the goal is to find the parameters θ that define the density function $f(\beta|\theta)$ where the functional form $f(\cdot)$ is given but parameters θ are unknown.

To approximate the probability value given a particular θ , Draw multiple β_r vectors from the distribution $f(\beta|\theta)$

$$P(Y_i = k) \approx \sum_{r=1}^R \frac{e^{\beta_r' x_{ik}}}{\sum_{l=1}^K e^{\beta_r' x_{il}}}$$

Plug this approximation into the log-likelihood Objective function and estimate θ value by doing optimization.

For mixed logit with repeated choices
 (panel data) $i = \text{individual}$ $t = \text{observation}$
 $k = \text{alternative}$

$$P(Y_{i1} = x_1, \dots, Y_{iT} = x_T) \\ = \int \prod_{t=1}^T \left(\frac{e^{\beta' x_{i,t}}}{\sum_{k=1}^K e^{\beta' x_{i,k,t}}} \right) f(\beta) d\beta$$

Panel data needs to account for the fact that the errors are correlated for the same individual over time.

Building more sophisticated models (even more)

```
M1 <- mlogit (Choice ~ CC + GN + NS + BU + FA + LD  
+ BZ + FC + FP + RP + PP + KA + SC + TS + NV + MA  
+ LB + AF + HU + PRICE - 1, data = S,  
upar = C(CC = 'n', GN = 'n', NS = 'n', BU = 'n',  
FA = 'n', LD = 'n', BZ = 'n', FC = 'n', FP = 'n',  
RP = 'n', PP = 'n', KA = 'n', SC = 'n', TS = 'n',  
NV = 'n', MA = 'n', LB = 'n', AF = 'n', HU = 'n',  
PRICE = 'n'), panel = TRUE, print.level =  
TRUE)
```

This fits a mixed logit model where the coefficients are treated as random variables (This is captured with upar (random parameters) argument where 'n' indicates it is modeled as a normal random variable; panel data captures the fact that we have multiple observations per individual.

Summary (M1)

This estimates for each random coefficient, a mean value and a standard deviation

Loglikelihood = -6531.3

$P1 \leftarrow \text{predict}(M1, \text{new data} = S)$

$\text{Predicted Choice 1} \leftarrow \text{apply}(P1, 1, \text{which.max})$

table(Predicted Choice 1, Actual)

	Ch 1	Ch 2	Ch 3	Ch 4
1	530	243	179	288
2	203	552	207	256
3	173	220	537	236
4	465	515	412	905

$$\text{Correct predictions} = \frac{530 + 552 + 537 + 905}{6000}$$

$$= \boxed{0.420}$$

Note that while the log-likelihood for the mixed logit model is better, in terms of predicting by simply looking at highest probability it is worse than MNL in this example.

$\text{Test Predict 1} \leftarrow \text{predict}(M1, \text{newdata} = \text{Test})$

$\text{Actual Choice 1} \leftarrow \text{subset}(\text{Safety}, \text{Task} > 12) [, "choice"]$

$\text{Predicted Choice 1} \leftarrow \text{apply}(\text{Test predict}, 1, \text{which.max})$

table(Predicted Choice 1, Actual Choice 1)

331	116	77	137
102	377	83	180
98	110	316	149
263	271	246	644

Correct predictions

$$= \frac{331 + 377 + 316 + 644}{3560}$$

$$= \boxed{0.4765}$$

The mixed logit model does a better job of predicting customers who are not interested in choosing any one of the offered options compared to MNL.

We can also compare out of sample (test) log-likelihood.

Prob Predict \leftarrow apply (TestPredict, 1, max)
MNL

Sum (log (Prob Predict))

-2703.78

Prob Predict 1 \leftarrow apply (TestPredict1, 1, max)
Mixed logit

Sum (log (Prob Predict 1))

-3086.894