

Test your knowledge of logit & big data analytics

1) `heat <- read.csv("Heating.csv")`

`head(heat)` This helps us visualise the first few observations in the dataset.

a) `install.packages("mlogit")` (If mlogit not installed before)
`library(mlogit)`

`data <- mlogit.data(heat, choice = "deprvar",
shape = "wide", varying = c(3:12))`

This creates a data object for mlogit to run where deprvar is the choice, wide indicates that each row is one choice situation & varying indicates the indexes of variables that are alternative specific

`model1 <- mlogit(deprvar ~ ic + oc - 1, data)`

Here ic & oc are used to predict choice where -1 indicates that no intercepts are needed

i) Both the coefficients of ic & oc are negative which makes sense since as the installation cost & operating cost for a system increases, the probability of choosing that system goes down.

ii) p value for $H_0: \beta_{ec} = 0$ is $< 2.2 \times 10^{-16}$

p value for $H_0: \beta_{oc} = 0$ is $< 2.2 \times 10^{-16}$

These p-values are very close to 0 indicating that we can reject the null hypothesis that the coefficients are 0.

iii) `pred1 <- predict(model1, newdata = data)`

This predicts the choice probabilities for each choice situation

`table (heat & depvar)/900`

ec	er	gc	gr	hp
0.071	0.093	0.636	0.143	0.055

These are the observed shares in the data set

`apply(pred1, 2, mean)`

This computes the average of the choice probabilities & gives:

ec	er	gc	gr	hp
0.104	0.051	0.516	0.24	0.087

While the model captures the data reasonably well, there are still differences in the results

$0.636/0.516$ $0.143/0.24$.

iv) $\text{Model 1} \$ \text{coef}["\text{oc}"] / \text{model 1} \$ \text{coef}["\text{ic}"]$

This computes the ratio of the $\beta_{oc} / \beta_{ic} = 0.739$

This implies that the decision makers are willing to pay \$0.73 higher in installation cost for 1 \$ reduction in operating cost. Note that it seems unreasonable for the decision-maker to pay only 73 cents for a one time payment for a 1 \$ reduction per year. economically.

v) $\text{data} \$ \text{lcc} \leftarrow \text{data} \$ \text{ic} + \text{data} \$ \text{oc} / 0.12$

This defines the new set of lifecycle costs in the data object for `mllogit`

$\text{model2} \leftarrow \text{mllogit}(\text{deprvar} \sim \text{lcc} - 1, \text{data})$

The log likelihood = -1248.7 & the variable `lcc` is significant.

c)

```
model3 <- mlogit( depvar ~ ic + oc, data,  
                   relevel = "hp" )
```

This forces hp to be the reference level &
the other alternative specific constants are relative
to this

```
pred3 <- predict(model3, newdata = data)
```

```
apply(pred3, 2, mean)
```

```
table(heat & depvar) / 900
```

In this case both match exactly

ec	er	gc	gr	hp
0.0711	0.0933	0.636	0.143	0.055

The presence of alternative specific constants
ensures that the average probabilities
equal the average share.

ii) $\text{model3} \$ \text{coef}["oc"] / \text{model3} \$ \text{coef}["ic"]$

In this case = 4.56 \$ for 1 \$ saving in
~~100~~ operating costs. This seems more reasonable
since operating cost is annual while installation
cost is one time.

iii) $\text{model4} \leftarrow \text{mlogit}(\text{depr} \sim ic + oc, \text{data},$
 $\text{reflevel} = "gr")$

Here gr is the reference level. Note that
in model 3, the intercept for gr is 0.308.
Hence in the new model, we would reduce
all the alternative specific constants downward
by 0.308. Note this does not change quality
of fit since probabilities are not modified
by adding a constant to all alternatives.

d)

$$i) \text{date\$iic} \leftarrow \text{date\$ic} / \text{date\$Income}$$

$$\text{models} \leftarrow \text{mlogit}(\text{depr} \sim \text{oc} + \text{iic}, \text{date})$$

Note that by `summary(models)`, we note that the new log likelihood value = -1010.2. In comparison to ~~-1008.2~~ -1008.2 is lower. This is worse in terms of fit since we want to maximize log likelihood. Also in the new model, installation cost divided by income is no longer significant while it was significant in the model in part c).

$$ii) \text{model 6} \leftarrow \text{mlogit}(\text{depr} \sim \text{oc} + \text{ic} | \text{Income}, \text{date})$$

Summary (model 6)

This also estimates a coefficient for each alternative that is income dependent. As income rises,

probability of choosing hot pump increases (zero coefficient)

relative to other. As income rises, probability

of choosing gas room drops relative to others. (most -ve)

None of these variables are significant however

e) $\text{heat}1 \leftarrow \text{heat}$

i) $\text{heat}1 \& \text{ic.hp} \leftarrow 0.9 \text{ heat}1 \& \text{ic.hp}$

$\text{data}1 \leftarrow \text{mlogit}.\text{data}(\text{heat}1, \text{choice} = "depvar",$

$\text{Shape} = "wide", \text{varying} = c(3:12))$

We will use model 3 to make predictions

$\text{pred}3a \leftarrow \text{predict}(\text{model}3, \text{newdata} = \text{data}1)$

$\text{apply}(\text{pred}3a, 2, \text{mean})$

<u>hp</u>	<u>ec</u>	<u>er</u>	<u>gc</u>	<u>gr</u>
0.064	0.704	0.092	0.63	0.142

These are the predicted shares under the new installation costs.

The market share of heat pumps goes to 0.0643 from 0.055 (6.43% from 5.5%)

ii)

$df \leftarrow \text{subset}(\text{heat}, \text{select} = c(3:12))$

Take only the installation & operating costs for the 5 original alternatives.

Define a new alternative with 2 new costs:

$$df\$ic.eci \leftarrow df\$ic.ec + 200$$

$$df\$oc.eci \leftarrow 0.75 \times df\$oc.ec$$

This new alternative (eci) has 200\$ more installation cost & 75% of the operating cost.

We will use model3 to estimate new choice probabilities.

$\text{model3}\$coef$ (Provides all coefficients)

We now compute new choice probabilities as follows:

$$df\$hpexp \leftarrow \exp(\text{model3}\$coef["oc"] \times df\$oc.hp + \text{model3}\$coef["ic"] \times df\$ic.hp)$$

$$df\$ecexp \leftarrow \exp(\text{model3}\$coef["oc"] \times df\$oc.ec + \text{model3}\$coef["ic"] \times df\$ic.ec + \text{model3}\$coef["ec:(intercept)"])$$

⋮

Repeat for all 6 alternatives

We now normalize by the sum of exponentials
to create the probabilities

```
df$sumexp <- apply(subset(df, select=c(13:17)), 1,  
                    sum)
```

```
df$sumnewexp <- apply(subset(df, select=c(13:18)),  
                      1, sum)
```

```
df$hp <- df$hpexp / df$sumexp (df$sumexp } 5  
:  
df$gr <- df$grexp / df$sumexp
```

```
df$hpnew <- df$hpexp / df$sumexpnewexp } 6  
df$ecinew <- df$eciexp / df$sumnewexp
```

```
df2 <- subset(df, select=c(26:31))
```

This contains new probabilities

```
marketsharenew <- apply(df2, 2, mean)
```

```
marketshareold <- apply(predict(model2, data), 2,  
                          mean)
```

New market shares

hp new	ec new	er new	gc new	gr new	ec new
0.049	0.063	0.083	0.57	0.128	0.103

Old market shares

hp	ec	er	gc	g
0.055	0.071	0.093	0.63	0.143

The most market share is drawn from gc from 0.63 to 0.57. This is the gas central system.

Note that from the independence of irrelevant alternatives property, the ratio of market shares remains the same irrespective of other alternatives in the set.

Note that in terms of % drop it draws roughly 10% from each system due to IIA. Note that you would assume the electric system should possibly draw more from electric systems rather than gas systems but this is not the case here.

2)

a) `electricity <- read.csv("Electricity.csv")`

`str(electricity)`

This dataframe has 4308 observations with 26 variables.

```
data1 <- mlogit.data(electricity, id.var = "id",  
  choice = "choice", varying = c(3:26),  
  shape = "wide", sep = "")
```

In this case there is one row for each choice situation.

```
model1 <- mlogit(choice ~ pf + d + loc + wk  
  + tod + ses - 1, data = data1,  
  vnames = c(pf = 'n', loc = 'n', wk = 'n',  
    tod = 'n', ses = 'n'), panel = TRUE)
```

We use the default settings to solve the mixed logit model with panel data to indicate multiple observations per individual.

`Summary(model1)`

i) The mean coefficient of contract length is around -0.18 indicating consumers prefer shorter contracts.
Mean price coefficient is -0.84. Therefore a

customer will pay $0.18 / 0.84 \approx 0.21$ cents per kWh to reduce contract length by 1 year.

The numbers could be slightly different for you

ii)

(cl)

The coefficient of length is normally distributed with mean $\underbrace{-0.18}_{cl}$ and standard deviation $\underbrace{0.31}_{sd.cl}$.

Share of people with negative

$$\text{coefficients} = P\left(\underbrace{\mu}_{-0.18} + \underbrace{\sigma}_{0.31} Z \leq 0\right)$$

$$= \text{pnorm}\left(-\frac{\text{model1\$coef}["cl"]}{\text{model1\$coef}["sd.cl"]}\right)$$

$$= 0.719$$

Roughly 72% of the population dislike long term contracts.

$$\text{bit}) \text{ pnorm}\left(-\text{model1\$coef}["pf"]/\text{model1\$coef}["sd.pf"]\right)$$

$$0.9999998 \quad (\text{very close to } 1)$$

Thus most of the population has negative price coefficients as should be expected.

model2 <- mlogit (choice ~ pf + cl + loc + wk +
tod + season - 1, data = data1, npar =
C(cl = 'n', loc = 'n', wk = 'n', tod = 'n',
seas = 'n'), panel = TRUE)

summary(model2)

The old model has log likelihood of -4089.6 while the new model has a log likelihood of -4110 (smaller as should be expected)

The estimated value of the price coefficient (pf) is -0.81 .

c) `model3 <- mlogit (choice ~ pf + cl + loc + wk + tod + seas - 1, data = data1, vpar = c(cl = 'n', loc = 'n', wk = 'u', tod = 'n', seas = 'n'), panel = TRUE)`

Summary(model3)

From the results we see that, the uniform distribution of wk is from 0.133 to 2.58 with mean = 1.36 .

The estimate price coefficient = -0.811

3)

a) Clearly for every value of k , the best subset selection method will have the smallest training sum of squared errors. This finds the global optimum set for each k .

b) This cannot be determined since a low training sum of squared errors does not necessarily translate to low test sum of squared errors.

c)

1) True. Since in forward stepwise selection at each step, we add only 1 variable to the previous set in an optimal manner.

2) True. Since you drop a variable from the set in backward stepwise selection

3) False

4) False

5) False

4)

a) `College <- read.csv("college.csv")`

`set.seed(1)`

`trainid <- sample(1:nrow(College), nrow(College)*0.8)`

`testid <- -trainid`

`train <- College[trainid,]`

`test <- College[testid,]`

`str(train)` $\left. \begin{array}{l} \text{ } \\ \text{ } \end{array} \right\} \begin{array}{l} 621 \text{ observations in} \\ \text{training set and 156 observation} \\ \text{in test set.} \end{array}$
`str(test)`

b) `model1 <- lm(Apps ~ ., data = train)`

`summary(model1)`

The total sum of squared error for the training set is obtained as

`sum(model1$residuals^2)`

The average sum of squared error is

`mean(model1$residuals^2) = 1061946.`

$\text{product1} \leftarrow \text{predict}(\text{model1}, \text{newdata} = \text{test})$

$\text{mse.test1} \leftarrow \text{mean}((\text{test}\$App - \text{product1})^2)$

This gives the test mean Squared Error

$= 1075064.$

c) library(leaps)

$\text{model2} \leftarrow \text{regsubsets}(App \sim ., \text{data} = \text{train},$

~~maxsteps~~ $\text{numexp} = 17,$

$\text{method} = "backward")$

$\text{Summary}(\text{model2})$

We see from the result that the subset of size 16 is obtained by dropping

P. Undergrad.

~~also observe that the best subset~~
~~is given by the variable R.~~

d) $\text{Summary}(\text{model2})\$adjr2$

$\text{plot}(\text{Summary}(\text{model2})\$adjr2)$

$\text{which.max}(\text{Summary}(\text{model2})\$adjr2)$ 12

$\text{Summary}(\text{model2})\$which$

In addition to the intercept, we would include Private Yes, Accept, Enroll, Top 10 perc, Top 25 perc, F. Undergrad, Outstate, Room. Board, PhD, Terminal, S.F. Ratio, Expend, Grad. Rate. We would drop P. Undergrad, Books, Terminal, perc. alumni, Personal.

e) $\text{model3} \sim \text{lm}(\text{Apps} \sim \text{Private} + \text{Accept} + \text{Enroll} + \text{Top10perc} + \text{Top25perc} + \text{F. Undergrad} + \text{Outstate} + \text{Room. Board} + \text{PhD} + \text{S.F. Ratio} + \text{Expend} + \text{Grad. Rate}, \text{data} = \text{train})$

Summary(model3)

This gives all the coefficients estimated values.

$\text{predict3} \leftarrow \text{predict}(\text{model3}, \text{newdata} = \text{test})$

$\text{mse3} \leftarrow \text{mean}((\text{test}\$Apps - \text{predict3})^2)$

This gives a test mean squared error = 1070293.

Since this error is less than 1075064, it does improve on the accuracy.

f) library(glmnet)

grid ← 10ⁿ seq(10, -2, length = 100)

X ← model.matrix(Apps ~ ., College)

Y ← College\$Apps

~~set.seed(1)~~

model4 ← glmnet(X[trainid,], Y[trainid],
lambda = grid)

plot(model4, xvar = "lambda")

g) set.seed(1)

cvmodel4 ← cv.glmnet(X[trainid,], Y[trainid],
nfolds = 10, lambda = grid)

cvmodel4\$lambda.min

This gives $\lambda = 0.497$

Signif(cvmodel4\$lambda, 4) Only plots to 4
significant digits

cvmodel4\$nonzero

There are 17 nonzero, complete model.

Here LASSO gives the full model.

Test error will be same as model 1.



