

Kumar, Sunil; Dabgotra, Apurba Vishal

## Article

# A latent class analysis on the usage of mobile phones among management students

Statistics in Transition New Series

## Provided in Cooperation with:

Polish Statistical Association

*Suggested Citation:* Kumar, Sunil; Dabgotra, Apurba Vishal (2021) : A latent class analysis on the usage of mobile phones among management students, Statistics in Transition New Series, ISSN 2450-0291, Exeley, New York, Vol. 22, Iss. 1, pp. 89-114, <https://doi.org/10.21307/stattrans-2021-005>

This Version is available at:

<https://hdl.handle.net/10419/236817>

## Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

## Terms of use:

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



<https://creativecommons.org/licenses/by-nc-nd/4.0/>

## A latent class analysis on the usage of mobile phones among management students

Sunil Kumar<sup>1</sup>, Apurba Vishal Dabgotra<sup>2</sup>

### ABSTRACT

In the past few years, wireless devices, including pocket PCs, pagers, mobile phones, etc., have gained popularity among a variety of users across the world and the use of mobile phones in particular, has increased significantly in many parts of the world, especially in India. Cell phones are now the most popular form of electronic communication and constitute an integral part of adolescents' daily lives, as is the case for the majority of mobile phone users. In fact, mobile phones have turned from a technological tool to a social tool. Therefore, the influence of cell phones on young people needs to be thoroughly examined. In this paper, we explore the attitude of young adults towards cell phones and identify the hidden classes of respondents according to the patterns of mobile phone use. The Latent Class Analysis (LCA) serves as a tool to detect any peculiarities, including those gender-based. LCA measures the value of an unknown latent variable on the basis of the respondents' answers to various indicator variables; for this reason, a proper selection of indicators is of great importance here. In this work, we propose a method of selecting the most useful variables for an LCA-based detection of group structures from within the examined data. We apply a greedy search algorithm, where during each phase the models are compared through an approximation to their Bayes factor. The method is applied in the process of selecting variables related to mobile phone usage which are most useful for the clustering of respondents into different classes. The findings demonstrate that young people display various feelings and attitudes toward cell phone usage.

**Key words:** backward greedy search algorithm (BGSA), latent class analysis (LCA), AIC, BIC.

### 1. Introduction

With the advancement in technology, the number of users of mobile phones have increased rapidly in the entire world. It has now become a crucial part of majority of lives of youth. This is because there are so many applications available on cell phones these days like internet access, sending e-mails, games, access to social networking sites like Facebook, listening to music, playing radio, reading books, dictionary and so on. Therefore, youth use their mobile phones more often for their free time. It was stated in the report of TRAI (Jan 2018), that India has about 1.012 billion mobile phone connections, which makes India's telecommunication market, the world's second largest in terms of number of

---

<sup>1</sup>Department of Statistics, University of Jammu, J&K, India. E-mail: sunilbhoulal06@gmail.com.  
ORCID: <http://orcid.org/0000-0003-0249-8415>.

<sup>2</sup>Department of Statistics, University of Jammu, J&K, India. E-mail: apurvavdabgotra@gmail.com.  
ORCID: <http://orcid.org/0000-0002-8056-7239>.

wireless connections after China. With youth population constituting half of the total population, India has become a fine breeding ground for highest cell connections. There exists a need to study and analyze the influence of cell phones on youth, because cell phones are responsible for modulating the thought process of any person especially that of the youngsters. In this paper, LCA is used for exploring the attitude of youth towards cell phones; to identify the hidden classes of respondents according to their mobile usage pattern and to arrive at the peculiarities in the cell phones usage, gender-wise, if any. The identification of most relevant variables related to mobile phone usage may aid in the evaluation of its user's attitude towards it and its impact on their lives in a way that usefully informs the role and vitality of cell phones in lives of today's youth.

LCA provides a tool for clustering and classification of individuals given the response pattern of a respondent to various questions or items of a questionnaire in a qualitative research. It is used for modelling and explaining the relationships between manifest or observed variable (may be dichotomous or polytomous) with respect to some unobserved or latent variables (may be dichotomous or polytomous) on the basis of data obtained in various kinds of surveys. LCA identifies unobservable (latent) subgroups within a population based on individuals' responses to different categorical observed variable. In addition to the above, LCA can also be used for the estimation of unknown parameters. The parameters of interests in any typical problem of latent class analysis are the unobserved proportion or size of the latent classes and the conditional item-response probabilities given the membership in a latent class. LCA will also provide estimate of the probabilities of respondents being misclassified by the questions. LCA was introduced in 1950 by Paul F. Lazarsfeld as a way of formulating latent attitudinal variables. Lazarsfeld performed LCA for building typologies (or clustering) based on dichotomous observed variables.

Then, Goodman (1974) proposed the estimation of LC model parameters using the maximum likelihood approach. Dempster et al. (1977) provided maximum likelihood estimation in the case of the observed and missing data involved in LCA. Haberman (1979) established the relationship between LC models and log-linear models for unknown cell counts frequency tables. Formann (1984, 1992) constructed a linear logistic LCA for dichotomous and polytomous variables. Mooijjaart (1992) presented the application of EM algorithm in LCA in a very detailed way. Vermunt (2010) proposed the inclusion of the covariates in the LC models, which make it possible to predict the LC membership probabilities by covariates through a logistic link. Biemer (2010) discussed the use of LCA as a survey error evaluation technique. Many theoretical developments and applications of LCA have been proposed in the recent past years, to encourage the research studies in LCA, example: Lanza (2013), Boduszek et. al (2014), Kumar (2015, 2016, 2017), Porcu (2017), Petersen (2019) and Sapounidis (2019).

Since LCA measures the value of an unknown latent variable given the responses made by the respondents to the various indicator variables, so proper selection of relevant indicator variables out of the set of all possible manifest variables is required in order to get efficient estimates of the unknown parameters. We have used, Raftery and Dean (2010) greedy search algorithm for the selection of indicator variables. This algorithm is used for checking single variable for inclusion into/exclusion from the set of selected clustering variables. The "greedy" search checks the inclusion of each single variable not currently

selected into the current set of selected clustering variables. The variable that has highest evidence of inclusion is proposed and, if its clustering evidence is stronger than the evidence against clustering it is included. At every exclusion step the “greedy” search option checks the exclusion of each single variable in the currently selected set of clustering variables and proposes the variable that has lowest evidence of clustering. The proposed variable is removed if its evidence of clustering is weaker than its evidence against clustering.

Through greedy search algorithm, we have identified the indicators variables which measures the latent variables. Following Baumgartner and Steenkamp (2006), the most promising way of accounting for extreme response bias is the inclusion of statistical techniques in the data analysis. Based on qualitative nature of the data it is not possible to use the classical least squares theory. LCA calculate conditional probabilities using Bayesian methodology.

To understand the application of methodological development, we have considered a study on the behaviour of mobile phone users (especially, young users). The paper is ordered as follows: Section 2 describes the framework of the greedy search algorithm followed by LCA model. Section 3 presents the data description and section 4 provides the analysis of mobile users data using the software application LCAvarsel (Fop and Murphy, 2017) and poLCA (Linzer and Lewis, 2011), the most complete and user-friendly package for the variable selection and the estimation of the latent class models. Finally, in section 5, we summarize our findings and recognize the probable areas of further research.

## 2. Methodology

### 2.1. Variable Selection

The problem of selecting the relevant clustering variables, in LCA can be modified as a model selection problem. Different models are specified by the role allotted to the variables in relevance to their relationship with the clustering variable  $X$ . Then, these models are compared by means of a model selection criterion and relevant clustering variables are, thus chosen accordingly in order to form the best model. The framework was introduced by the work of Raftery and Dean (2006). The authors propose a procedure where,  $Y$  (the set of all manifest variables) is partitioned into the following subsets of variables:

1.  $Y^C$ , the set of current clustering variables;
2.  $Y^P$ , the variable(s) proposed to be added or removed from the set of clustering variables;
3.  $Y^{NC}$ , the set of other variables not relevant for clustering.

Given this partition of  $Y$  and the (unknown) clustering membership  $X$ , we can modify the question of effectiveness of  $Y^C$  for clustering as a question of model selection. The question, thus becomes to choose one of two different models, i.e.  $M_A$  which assumes that  $Y^P$  is useful for clustering and  $M_B$  which assumes that it is not. In model  $M_B$ , the set of variable(s) proposed to be added or removed from the set of clustering variables  $Y^P$  is conditionally independent of the cluster memberships  $X$  given the variables  $Y^C$  is already in the model. In model  $M_A$ , this is not the case. In both models, the set of other variables

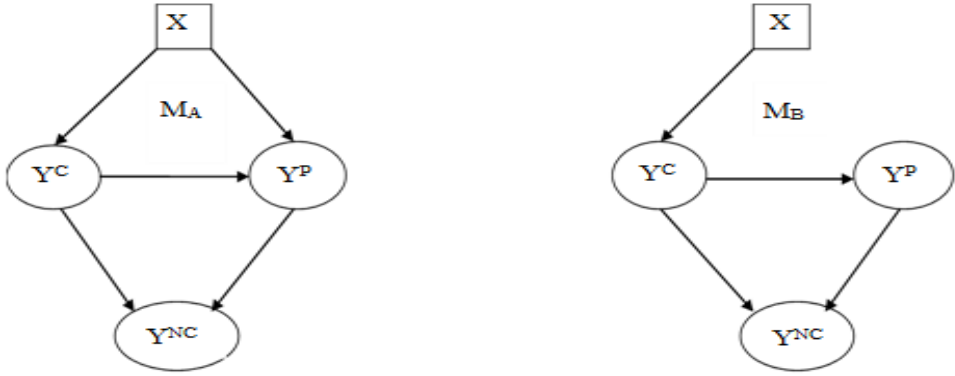
considered  $Y^{NC}$  is conditionally independent of cluster membership  $X$  given  $Y^C$  and  $Y^P$ , but may be associated with  $Y^C$  and  $Y^P$ . Then, the decision for inclusion or exclusion of the variable belonging to  $Y^P$  is taken by comparing the models (Figure 1):

$$\begin{aligned} M_A : P(Y|X) &= P(Y^C, Y^P, Y^{NC}|X) \\ &= P(Y^C, Y^P|X) P(Y^{NC}|Y^C, Y^P), \end{aligned} \quad (1)$$

$$\begin{aligned} M_B : P(Y|X) &= P(Y^C, Y^P, Y^{NC}|X) \\ &= P(Y^C|X) P(Y^P|Y^C) P(Y^{NC}|Y^C, Y^P), \end{aligned} \quad (2)$$

where,  $X$  is the (unobserved) set of cluster memberships. Model  $M_A$  implies that  $Y^P$  does provide information about clustering membership, beyond that given just by  $Y^C$ . Model  $M_B$  specifies that, given  $Y^C$ ,  $Y^P$  is independent of the cluster memberships (defined by the unobserved variables  $X$ ), i.e.,  $Y^P$  gives no further information about the clustering. In model  $M_A$ ,  $Y^P$  is useful for clustering and the joint distribution  $P(Y^C, Y^P|X)$  corresponds to a Gaussian mixture distribution; on the other hand,  $M_B$  states that  $Y^P$  does not depend on the clustering variable  $X$  and the conditional distribution  $P(Y^P|Y^C)$  corresponds to a linear regression.

**Figure 1:** Graphical Representation of Models  $M_A$  and  $M_B$  for Variable Selection.



An important feature of the framework formulation is that in  $M_B$  the irrelevant variables are not required to be independent of the clustering variables. This criterion allows to discard redundant variables related to the clustering ones but not to the clustering itself. Models  $M_A$  and  $M_B$  are compared via an approximation to the Bayes factor. The Bayes factor,  $B_{AB}$ , for  $M_A$  against  $M_B$  based on the data  $Y$  is given by

$$B_{AB} = \frac{P(Y|M_A)}{P(Y|M_B)}, \quad (3)$$

where,  $P(Y|M_i)$  is the integrated likelihood of model  $M_i$  ( $i=A,B$ ), namely

$$P(Y|M_i) = \int P(Y|\Theta_i, M_i) P(\Theta_i|M_i) d\Theta_i,$$

where ,  $\Theta_i$  is the vector-valued parameter of model  $M_i$ , and  $P(\Theta_k|M_i)$  is its prior distribution (Kass and Raftery, 1995).

Let us now consider the integrated likelihood of model  $M_A$ ,  $P(Y|M_A) = P(Y^C, Y^P, Y^{NC}|M_A)$ . From (1), the model  $M_A$  is specified by two probability distributions: the latent class model that specifies  $P(Y^C, Y^P|\Theta_1, M_A)$ , and the distribution  $P(Y^{NC}|Y^C, Y^P, \Theta_1, M_A)$ . We denote the parameter vectors that specify these two probability distributions by  $\Theta_{11}$ ,  $\Theta_{12}$ , and we assume that their prior distributions are independent. Then, the integrated likelihood factors as follows:

$$\begin{aligned} P(Y|M_A) &= P(Y^C, Y^P, Y^{NC}|M_A) \\ P(Y|M_A) &= P(Y^C, Y^P|M_A)P(Y^{NC}|Y^C, Y^P, M_A), \end{aligned}$$

$$\begin{aligned} \text{where , } P(Y^{NC}|Y^P, Y^C, M_A) &= \int P(Y^{NC}|Y^P, Y^C, \Theta_{12}, M_A)P(\Theta_{12}|M_A)d\Theta_{12} \\ \text{and } P(Y^C, Y^P|M_A) &= \int P(Y^C, Y^P|\Theta_{11}, M_A)P(\Theta_{11}|M_A)d\Theta_{11}. \end{aligned}$$

Similarly, we obtain

$$P(Y|M_B) = P(Y^C|M_B)P(Y^P|Y^C, M_B)P(Y^{NC}|Y^C, Y^P, M_B).$$

The prior distribution of the parameter,  $\Theta_{12}$ , is assumed to be the same under  $M_A$  as under  $M_B$ . It follows that

$$P(Y^{NC}|Y^P, Y^C, M_A) = P(Y^{NC}|Y^P, Y^C, M_B).$$

We thus have

$$B_{AB} = \frac{P(Y^C, Y^P|M_A)}{P(Y^C|M_B)P(Y^P|Y^C, M_B)} \quad (4)$$

which has been greatly simplified by the cancellation of the factors involving the potentially high-dimensional  $Y^{NC}$ . The integrated likelihoods in (4) are still hard to evaluate analytically, so we approximate them using the BIC approximation given as:

$$BIC(Y|M_i) = 2 \times \log(\max. \ln L) - (p) \times \log(n)$$

where , i= A, B

max. lnL : maximum log- likelihood

p : number of parameters

n : number of observations

The above equation provides the BIC approximates for model  $M_A$  & model  $M_B$  and thus leading to the following criterion:

$$BIC_{diff} = BIC_A - BIC_B.$$

The clustering variables are selected using a stepwise algorithm of greedy search algorithm. At each stage it searches for the variable to add that most improves the clustering as measured by BIC, and then assesses whether one of the current clustering variables can be dropped. At each stage, the best combination of the number of groups and clustering model is chosen. For performing greedy search algorithm, we first have to choose the value of no. of classes for each of the latent variable,  $r=1,2,...,R$ ; so that , our model is identifiable. For a given number of variables, not all the models specified by assigning different values to  $r$

are identifiable. In fact, a necessary (though not sufficient) condition to the identifiability of a model with  $r$  latent classes is

$$\prod_{m=1}^M C_m > \left( \sum_{m=1}^M C_m - M + 1 \right) r, \quad (5)$$

with  $C_m$  the number of categories taken by variable  $X_m$  [Goodman (1974)]. Thus when selecting the number of classes, hereafter we will consider values of  $r$  for which this identifiability condition holds.

Set  $R$  (max.  $r$ ), the maximum number of clusters to be considered for the data. Make sure that this number is identifiable for the data, i.e.  $R$ , the maximum number of latent classes should satisfy the identifiability condition in (5) for the set of variables currently taken into consideration in fitting the LCA model. If a latent class model on the set of all variables is identifiable for  $R > 1$ , we then estimate the model. For each category of each variable, we calculate the variance of its probability across groups. Then, for each variable, we add up these variances and rank the variables according to that sum. The reason behind this, is that the variables with high values of this sum have high between-group variation in probability, and hence they may be more useful for clustering. Given this ranking we choose the top  $k^*$  variables, where  $k^*$  is the smallest number of variables that allow a latent class model with  $R > 1$  to be identified. This will provide the starting  $Y^C$ . The other variables can be left in their ordering based on variability for future order of their inclusion in the algorithm. In fitting the LCA model we perform multiple runs with random starting values. Also, in this case the aim is to allow the search for the global maximum of the log-likelihood rather than a local one; then the model with the greatest log-likelihood is retained. The following are the inclusion and exclusion steps of a new variable:

**Inclusion Step:** We look at each variable in  $Y^{NC}$  as to be a new variable under consideration for inclusion into  $Y^P$  and the variable that has the highest evidence of inclusion is proposed. Then, calculate the difference in BIC for models  $M_A$  and  $M_B$  given the current  $Y^C$ . If the variable's BIC difference is:

- below 0, do not include the variable in  $Y^C$  and remove variable from  $Y^{NC}$ ;
- above 0, include variable in  $Y^C$  and stop inclusion step.

If we reach the end of the list of variables in  $Y^{NC}$ , the inclusion step is stopped.

**Exclusion Step:** We look at each variable in  $Y^C$  as to be a new variable under consideration for inclusion into  $Y^P$  (with the remaining variables in  $Y^C$  not including current  $Y^P$  now defined as  $Y^C$  in  $M_A$  and  $M_B$ ) and the variable that has the lowest evidence of being in the set is proposed. Calculate the difference in BIC for models  $M_A$  and  $M_B$ . If the variable's BIC difference is:

- below 0, remove the variable from (the original)  $Y^C$  and place it under  $Y^{NC}$  and stop the exclusion step;
- above 0, do not remove the variable from (the original)  $Y^C$ .

If we reach the end of the list of variables in  $Y^{NC}$  the exclusion step is stopped.

If  $Y^C$  remains the same after consecutive inclusion and exclusion steps the greedy search algorithm stops because it has converged.

## 2.2. Latent Class Analysis

After identifying the variables for each latent variables, we next perform LCA for exploring the attitude of youth towards cell phones and to categories the respondents according to the pattern of their mobile phone usage. LCA has been used as a multivariate statistical tool for the study. LCA models comprises two types of probabilities which include

- the probability indicating the likelihood of a response by respondents in each of the classes and
- the probability representing the latent class size or the proportion of individuals who are members of a particular latent class.

The former one represents the probability of a particular responses to a manifest variable, conditioned on latent class membership and can be interpreted as factor loading for Factor Analysis, in which both the observed or latent variables are continuous. LCA provides a clustering of individuals in a population, based on the response patterns of individuals to the different observed variables.

The assumptions of the standard LC model (Biemer,2010) are as follows:

1. The sample can be treated as if it was a simple random sample without replacement from an infinite population, i.e. data is sampled without replacement from a large population units using SRS (Simple Random Sampling).
2. The indicators are locally independent within a latent class, means all the indicator variables have nothing in common except latent variable, i.e. after accounting for latent variable  $X$ , there is no association between indicator variables.
3. The response probabilities are homogeneous, i.e. the probabilities of selecting any two units (individuals) from the population are same.
4. The indicator variables are univocal, i.e. the indicator variables can measure one and only one latent variable.

Following the notation used by Linzer and Lewis (2011), suppose we have  $J$  polytomous categorical manifest variables ( the observed variable) each of which contain  $K_j$  possible outcomes, for individuals  $i = 1, 2, 3, \dots, N$ . Let  $Y_{ijk}$  be the observed values of the  $J$  manifest variables such that

$$\left\{ \begin{array}{ll} Y_{ijk} = 1 : & \text{if } i^{th} \text{ respondent give the } k^{th} \text{ response to the } j^{th} \text{ variable} \\ Y_{ijk} = 0 : & \text{otherwise} \end{array} \right\}$$

where,  $j=1,2,\dots,J$  and  $k=1,2,\dots,K_j$ .

The LC models approximates the observed joint distribution of the manifest variables as the weighted sum of a finite number,  $R$ , of constituent cross-classification tables. Let  $\pi_{jrk}$  denote the cross-conditional probability that an observation in class  $r=1,2,\dots,R$  produces the  $k^{th}$  outcome on the  $j^{th}$  variable with

$$\sum_{k=1}^{K_j} \pi_{jrk} = 1.$$

Let  $p_r$  be the prior probabilities of latent class membership, as they represent the unconditional probability that an individual will belong to each class before taking into account



the responses  $Y_{ijk}$  provided on the manifest variables. The probability that an individual  $i$  in class  $r$  produces a particular set of  $J$  outcomes on the manifest variables, assuming conditional independence of the outcomes  $Y$  given class membership, is the product

$$f(Y_i; \pi_r) = \prod_{j=1}^J \prod_{k=1}^{K_j} (\pi_{jrk})^{Y_{ijk}}, \quad (6)$$

The probability density function across all classes is the weighted sum

$$f(Y_i | \pi, p) = \sum_{r=1}^R f(Y_i; \pi_r) = \sum_{r=1}^R P_r \prod_{j=1}^J \prod_{k=1}^{K_j} (\pi_{jrk})^{Y_{ijk}}. \quad (7)$$

The parameters  $P_r$  and  $\pi_{jrk}$  are estimated by the latent class model.

The unknown parameters of the LC models can be estimated by maximizing the Log likelihood function with respect to  $p_r$  and  $\pi_{jrk}$ , using the expectation-maximization (EM) algorithm (Dempster et al. (1977), McLachlan and Peel (2000) and Linzer and Lewis (2011)) some other algorithms can also be considered like a Newton-Raphson algorithm or a hybrid form of these two algorithms (McLachlan and Krishnan, 2008). In any case the algorithm is initialized through a set of randomly generated starting values and there is no guarantee of reaching the global maximum. For this reason, it is usually a good practice to run the procedure a number of times and select the best solution [Bartholomew et al. (2011)]. Given estimates  $\hat{P}_r$  and  $\hat{\pi}_{jrk}$  of  $P_r$  and  $\pi_{jrk}$  respectively, the posterior probability that each individual belongs to each class, conditional on the observed values of the manifest variables, is calculated by

$$\hat{P}(r_i | Y_i) = \frac{\hat{P}_r f(Y_i; \hat{\pi}_r)}{\sum_{q=1}^R \hat{P}_q f(Y_i; \hat{\pi}_q)}, \quad (8)$$

where,  $r_i \in (1, 2, \dots, R)$ .

It is important that the condition  $R \sum_j (K_j - 1) + (R - 1) \leq n$  on the number of parameters should hold. Also,  $R \sum_j (K_j - 1) + (R - 1) \leq (3^{10} - 1)$ , i.e. one fewer than the total number of cells in the cross-classification table of the manifest variables, as then the latent class model will be unidentified. Under the assumptions of multinomial distribution, the log likelihood function can be given as

$$\ln L = \sum_{i=1}^n \ln \sum_{r=1}^R p_r \prod_{j=1}^J \prod_{k=1}^{K_j} (\pi_{jrk})^{Y_{ijk}}. \quad (9)$$

LCA not only builds a measurement or classification model but it also explains a relation of the class membership to explanatory variables by including covariates (single grouping variable or a combination of grouping variables) (Vermunt, 2010) in the model. These explanatory variables are referred to as covariates, predictors, external variables, independent variables, or concomitant variables. In a more explanatory study, one may wish to build a predictive or structural model for class membership whereas in a more descriptive study the aim would be to simply profile the latent classes by investigating their association with

external variables. Grouping variables can be used in LC models in order to model the unexplained heterogeneity in the data. In that case latent class membership probabilities are predicted by covariates through a logistic link.

Once the LC models are built then the next step is to select the optimum model as different LC models have a different number of latent classes. Usually, models with more parameters (i.e more latent classes) provide a better fit, and more parsimonious models tend to have a somewhat poorer fit. So, there is always very close agreement between goodness of fit and parsimony of the latent class models. We can test the goodness of fit of an estimated LCA models by the Pearson Chi-square ( $\chi^2$ ) or the Likelihood Ratio Chi-square ( $L^2$ ). However, the likelihood ratio Chi-square test, although extensively used in the statistical literature, has a number of important limitations. These limitations can be controlled by making use of several information criteria, such as the Akaike information criterion (AIC) (Akaike (1973)) and Bayesian information criterion (BIC) (Schwartz (1978)), each of which is designed to penalize models with larger numbers of parameters. AIC and BIC on the number of parameters in the model:

$$AIC = L^2 - 2 \times d.f. \quad \text{and} \quad BIC = L^2 - d.f. \times \ln(n),$$

where,  $n$  is the sample size.

These information criteria are commonly used for selecting the optimal number of latent classes in a model. By comparing models with a different number of latent classes, a model with lower AIC and BIC is selected.

### 2.3. Latent Regression Models

The latent class regression model generalizes the basic latent class model by permitting the inclusion of covariates to predict individuals' latent class membership (Dayton and Macready, 1988; Hagenaars and McCutcheon, 2002). This is the so-called "one-step" technique for estimating the effects of covariates, because the coefficients on the covariates are estimated simultaneously as part of the latent class model. Covariates are included in the latent class regression model through their effects on the priors  $p_r$ . In the basic latent class model, it is assumed that every individual has the same prior probabilities of latent class membership. The latent class regression model, in contrast, allows individuals' priors to vary depending upon their observed covariates.

Denote the mixing proportions in the latent class regression model as  $p_{ri}$  to reflect the fact that these priors are now free to vary by individual. It is still the case that  $\sum_r p_{ri} = 1$  for each individual. To accommodate this constraint, a generalized (multinomial) logit link function can be employed for obtaining the effects of the covariates on the priors (Agresti, 2002). Let  $X_i$  represent the observed covariates for individual  $i$ . First latent class is arbitrarily selected as a "reference" class and assumes that the log-odds of the latent class membership priors with respect to that class are linear functions of the covariates. Let  $\beta_r$  denote the vector of coefficients corresponding to the  $r^{th}$  latent class. With  $S$  covariates, the  $\beta_r$  have length  $S + 1$ ; this is one coefficient on each of the covariates plus a constant. Because the first class is used as the reference,  $\beta_r = 0$  is fixed by definition.

Then,

$$\ln(p_{2i}/p_{1i}) = X_i\beta_2$$

$$\ln(p_{3i}/p_{1i}) = X_i\beta_3$$

.

.

.

$$\ln(p_{Ri}/p_{1i}) = X_i\beta_R$$

Following some simple algebra, this produces the general result that

$$p_{ri} = p_r(X_i; \beta) = \frac{e^{X_i\beta_r}}{\sum_{q=1}^R e^{X_i\beta_q}}. \quad (10)$$

The parameters estimated by the latent class regression model are the  $R - 1$  vectors of coefficients  $\beta_r$  and as in the basic latent class model, the class-conditional outcome probabilities  $\pi_{jrk}$ . Given estimates  $\hat{\beta}_r$  and  $\hat{\pi}_{jrk}$  of these parameters, the posterior class membership probabilities in the latent class regression model are obtained by replacing the  $p_r$  in Eq. 8 with the function  $p_r(X_i; \beta)$  in Eq. 10:

$$\hat{P}(r|X_i; Y_i) = \frac{p_r(X_i; \hat{\beta})f(Y_i; \hat{\pi}_r)}{\sum_{q=1}^R p_q(X_i; \hat{\beta})f(Y_i; \hat{\pi}_q)}. \quad (11)$$

The number of parameters estimated by the latent class regression model is equal to  $\sum_j^R (K_j - 1) + (S + 1)(R - 1)$ . The same considerations mentioned earlier regarding model identifiability also apply here. The parameters of the LC regression model are estimated in the same way as simple LC models, the only difference is that  $p_r$  is replaced by  $p_r(X_i; \beta)$ .

### 3. Data Description

The sample included management students at post-graduation level who were using smartphones for more than one year. The methodology evolved in the research was divided into three phases. The first phase was about understanding the role of smartphone in students' life with the identification of dimensions that influence their behaviour on day-to-day basis. That led to the development of an instrument based on the literature review of studies carried out in the past. The self-report instrument from this development is supplied to the students that makes the second stage of the study, i.e. data collection. The instrument was created with the help of Google form and thus, online data was collected. The survey was administered in a top-rank state University of Jammu, India, which boasts about its excellent technological infrastructure. The third and final phase was editing, coding and analysis of data. A total of 214 responses were used in analysis. Male respondents were 71.5% and female respondents were 28.5%. The analysis was carried out with statistical software SPSS-25 and R software.

The questionnaire contained a total of 28 seven-point Likert type scale item ranging from 'completely unimportant' to 'completely important'; 2 three-point Likert type scale items ranging from 'unimportant' to 'important' and 1 two-point scaling. The questionnaire also consists of Qualitative questions on Usage of cell phones; Necessity of cell phones in modern time; Cost efficiency of cell phones over landlines; Safety reasons for carrying cell phones; Reliance on the cell phones; Dependency on the cell phones; Non calling functionality of cell phones.

Variables on mobile phone usage behaviour viz. information access, personal safety, financial incentives, social interactions, parental contacts, time management, dependency, reputations, gender, brand importance, etc, are taken to be the observed or manifest variable for the analysis. All the manifest variables have polytomous (i.e. 7) response options except for 1 that have binary responses. Accordingly, 7 latent variables have been considered on the basis of different manifest (observed) variables. A detailed description of the variables used in this paper is in Appendix A.

A pilot study was also conducted, in which the reliability and validity of the questionnaire is evaluated using SPSS software and it shows an internal consistency reliability of about 0.89 (i.e. Cronbac alpha = 0.89), which means that the internal consistency of the questionnaire is good. Also, the construct validity of the questionnaire is also quite high.

## 4. Results

Our data set consists of 30 variables related to the different behaviour of youths over mobile usage and 1 variable related to their gender. Due to the different behaviour our population consists of highly heterogeneous variable, which results in the violation of LCA assumption of identifiability. Thus, we formulated a different combination of variables in the form of 7 subsets on the basis of correlation. Thus, there may be 7 latent variables and the variable selection procedure has been performed for each of them, so as to get the relevant set of indicator variables for each of them. For selecting variables, we have performed analysis by using backward greedy search algorithm (BGSA). BGSA is performed in the following way:

First, we find the number of latent classes to be considered for the data, in order to make LCA models identifiable. Because, at each step, the BGSA considers only latent class analysis models for which the identifiability condition described in equation (5) holds. Then, we choose the maximum possible number of latent classes for which the LC model is identifiable and perform BGSA using LCAvarsel package of R software. In this case the value of upper as well as lower is 0 by default. The individual step results of the variable selection procedure starting with a relevant set of variables are given in Table 1.

The greedy search algorithm starts with two successive removal steps, then it iterates alternating between removal and inclusion step. It stops when the set of relevant predictors remains unchanged after consecutive removal and inclusion steps. Column 2 of Table 1 shows the possible set of variables proposed for each of the latent variable in the greedy search algorithm. For latent variable  $X_1$ , initially  $a_7$  (having highest BIC) is proposed for the removal and then the difference between the BIC of the two models is computed.

**Table 1:** Stepwise result of variable selection algorithm

Latent variables	Proposed variables	Step	Variable	BIC diff.	Decision	Selected variables
$X_1$	$a_1, a_2, a_4, a_7, g$	Remove	$a_7$	190.63	Accepted	$a_1, a_2, a_4, g$
		Remove	NA	NA	NA	
		Add	$a_7$	-135.90	Rejected	
$X_2$	$a_3, a_5, a_7, g$	Remove	$g$	-185.33	Rejected	$a_3, a_5, a_7, g$
$X_3$	$a_8, a_9, a_{10}, a_{11}, a_{12}, g$	Remove	$a_8$	115.34	Accepted	$a_9, a_{10}, a_{11}, a_{12}, g$
		Remove	$a_{10}$	-31.69	Rejected	
		Add	$a_8$	-148.60	Rejected	
$X_4$	$a_{13}, a_{14}, a_{15}, a_{16}, a_{17}, a_{18}, a_{19}, a_{26}, g$	Remove	$a_{26}$	9.13	Accepted	$a_{13}, a_{14}, a_{15}, a_{16}, a_{17}, a_{18}, a_{19}, g$
		Remove	$g$	-6.98	Rejected	
		Add	$a_{26}$	-9.13	Rejected	
$X_5$	$a_{20}, a_{21}, a_{22}, a_{23}, g$	Remove	$a_{23}$	91.94	Accepted	$a_{20}, a_{21}, a_{22}, a_{23}, g$
		Remove	NA	NA	NA	
		Add	$a_{23}$	7.79	Accepted	
		Remove	NA	NA	NA	
$X_6$	$a_{24}, a_{25}, a_{26}, a_{27}, g$	Remove	$a_{26}$	2.41	Accepted	$a_{24}, a_{25}, a_{27}, g$
		Add	$a_{26}$	-86.82	Rejected	
$X_7$	$a_{29}, a_{30}, g$	Remove	$a_{29}$	52.68	Accepted	$a_{29}, a_{30}, g$
		Add	$a_{29}$	7.52	Accepted	

For the removal step the LCAversal computes the difference  $BIC_B - BIC_A$ , that is why its decision rule's inequalities are reversed, i.e. if the BIC difference comes out to be greater than 0, then we remove the variable. For latent variable  $X_1$ , it is envisaged from Table 1, that removal of  $a_7$  is accepted as its  $BICdiff. > 0$  and the variable  $a_7$  is removed from the set of clustering variables  $Y^C$  and placed in  $Y^{NC}$ . When performing the stepwise selection, for some combinations of clustering variables and the number of classes, it could happen that a step of the variable selection procedure could not be performed because no latent class model is identifiable on any of the possible clustering sets. In such case, the step is not performed and a NA is returned. This is what happened at the next iteration of variable selection for  $X_1$ .

After two consecutive removal steps the inclusion step is performed and the removed variable  $a_7$  is again proposed for the inclusion in  $Y^C$ , as it is the only variable in  $Y^{NC}$ . Then again the difference between BIC of the two models is computed and if it comes out to be greater than 0, we have to include the variable. But in this case, the value of  $BICdiff. < 0$ , so the decision of adding  $a_7$  in  $Y^C$  is rejected. The selected indicator variables corresponding to each latent variable are shown in Table 1. Next, we perform Latent Class Analysis on the selected indicator variables for each latent variable in order to study the estimated heterogeneity of the population along with the number of latent classes.

From the results of BGSA, we have 7 latent variables, namely Work efficiency, Up to date, superiority over landlines, Safety/Security, Dependency, Negatives, Functionality; with selected, indicators variables. Further, we perform LCA on 7 latent variables in order to explain the heterogeneity among the sub-populations and conditional item-response probabilities, using polCA (Linzer and Lewis, 2011) package of R. The selection of an

appropriate number of latent classes for each of the latent variable can be made by comparing the values of BIC or AIC, a suitable measure for variable selection. Models considered in our analysis are simple extensions of the basic LC models as proposed by Biemer and Wiesen (2002) with grouping variables. The use of grouping variables has been suggested by Hui and Walter (1980) to either ensure that the model is identifiable, or to improve the fit of the model or to reduce the effect of unobserved heterogeneity. Following Hui-Walter approach, we choose gender, brand and addiction as grouping variables. Table 2 shows the goodness of fit statistics and other statistics with optimal number of latent classes for each latent variable.

From Table 2 it is clear that the data set is best fitted for a 3-class model for all the latent variables except for  $X_7$ , which has 2-class best fitted model, as the corresponding BIC as well as AIC values for each of the latent variables are the lowest. The data were inconsistent with several other possible models for each of the latent variable.

**Table 2:** The Model diagnostics for different latent variables

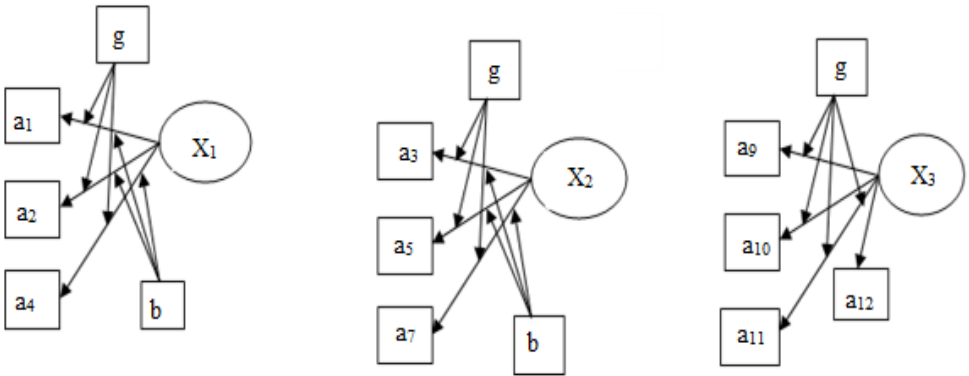
Latent variables	No. of opt. classes	Grouping variables	Residual d.f.	Max. log-likelihood	$\chi^2$	AIC	BIC
$X_1$	3	-	158	-1010.415	310.0607	2134.849	2355.075
		g	156	-1006.924	311.4976	2132.83	2349.325
		g&b	152	-1003.858	336.3663	2131.716	2340.407
$X_2$	3	-	158	-1072.475	267.0177	2279.204	2482.431
		g	156	-1068.602	279.4936	2266.949	2475.444
		g&b	152	-1064.097	312.6338	2252.195	2460.885
$X_3$	3	-	140	-1315.987	2585.898	2819.974	3169.056
		g	138	-1293.091	2695.378	2738.182	2993.996
		-	86	-1807.442	696845.5	3870.883	4301.728
$X_4$	3	g	84	-1801.077	1006308	3862.154	4299.73
		-	140	-1316.358	1023.609	2780.716	3029.798
		g	138	-1310.409	1559.712	2772.819	3028.633
$X_5$	3	g&a	134	-1300.755	1726.56	2761.51	3020.788
		-	158	-1148.867	285.8787	2369.734	2578.228
		g	156	-1128.732	286.4188	2345.464	2558.691
$X_6$	3	g&a	152	-1105.955	329.6206	2335.911	2544.601
		g	8	-781.8273	50.3265	1593.655	1698.294
		g&b	4	-766.1268	45.0985	1585.254	1688.356
$X_7$	2	g&a	4	-766.1388	51.1466	1585.278	1688.38
		g&b&a	16	-757.6595	94.3307	1579.319	1687.03

*d.f.* : degrees of freedom  
 $\chi^2$  : Chi – square value  
*AIC* : Akaike information criterion  
*BIC* : Bayesian information criterion

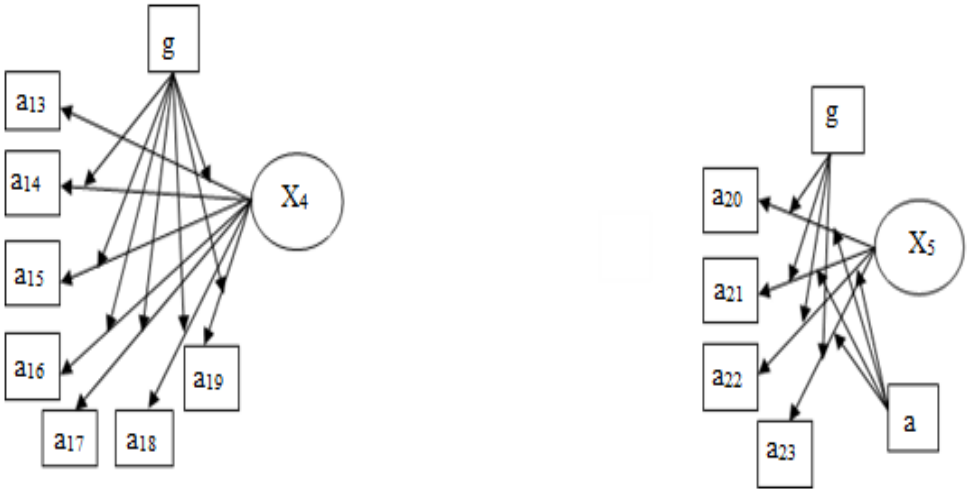
Since, we have heterogeneous indicator variables measuring the different latent variables and different grouping variables, for this reason various models were considered and

tested for their identifiability. Model diagnostics are provided in tTable 2 and the path models of identifiable and efficient models are provided in Figure 2, 3 and 4, for each of the latent variable.

**Figure 2:** Path models for latent variables  $X_1, X_2$  and  $X_3$

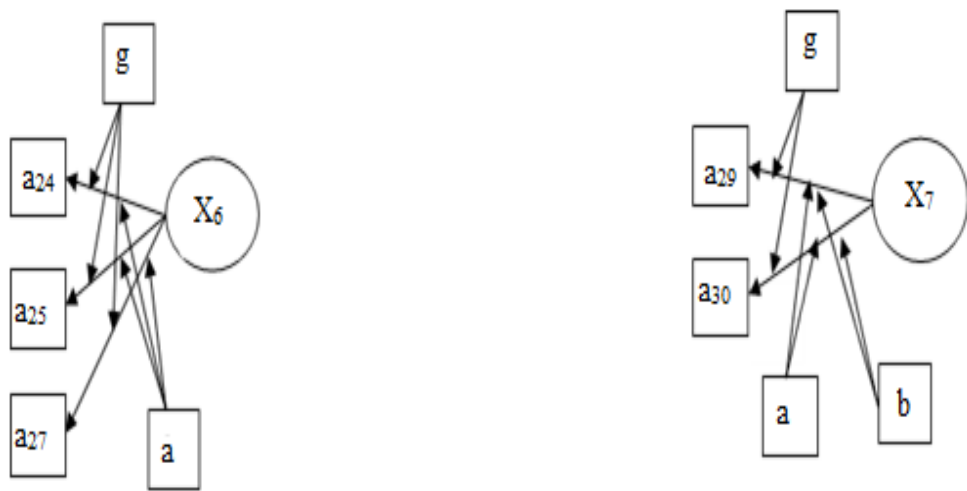


**Figure 3:** Path models for latent variables  $X_4$  and  $X_5$



From the results of Table 2, it is clear that the data set is best fitted for a 3-class model for all the latent variables except for  $X_7$ , which has 2-class for the best fitted model. Therefore, the underlying latent classes can be identified as “improve“(class 1), “same” (class 2) and “worsen” (class 3) for the latent variables  $X_1, X_2, X_3, X_4, X_5$  and  $X_6$ . And for latent variable  $X_7$ , the underlying latent classes can be identified as “non-calling” (class1) and “calling” (class2).

**Figure 4:** Path models for latent variables  $X_6$  and  $X_7$



**Table 3:** Estimated class membership for different latent variables

Latent variables	Grouping variables	Latent class 1	Latent class 2	Latent class 3
Work efficiency ( $X_1$ )	-	0.1946	0.2987	0.5067
	g	0.2846	0.5149	0.2005
	g&b	0.4871	0.3453	0.1676
Up to date ( $X_2$ )	-	0.2514	0.4393	0.3093
	g	0.1472	0.5115	0.3785
	g&b	0.3724	0.1721	0.4555
Superiority over landlines ( $X_3$ )	-	0.3883	0.289	0.3227
	g	0.334	0.2344	0.4316
Safety/security ( $X_4$ )	-	0.1159	0.5639	0.3202
	g	0.1605	0.1577	0.6818
	-	0.3645	0.1916	0.4439
Dependency ( $X_5$ )	g	0.1663	0.5233	0.3104
	g&a	0.2919	0.2376	0.4705
	-	0.4879	0.1743	0.3379
Negatives ( $X_6$ )	g	0.1755	0.4936	0.3309
	g&a	0.3034	0.3236	0.373
	-	0.374	0.626	-
Functionality ( $X_7$ )	g	0.1962	0.4416	0.3622
	g&b&a	0.2828	0.7172	-

The estimated class membership probabilities from the latent class model are summarized in Table 3. For latent variable  $X_1$ , the estimated class membership probabilities with grouping variable gender are 0.2846 for class 1, 0.5149 for class 2 and 0.2005 for class 3, while the estimated class membership probabilities of the same with the inclusion of 1 more





	Indic- tors	Latent classes	Categories						
			1	2	3	4	5	6	7
X <sub>4</sub>	a <sub>11</sub>	1	0.0551	0.1119	0.1106	0.0840	0.2054	0.3647	0.0682
		2	0.0000	0.0000	0.0218	0.4234	0.2257	0.0000	0.3292
		3	0.0115	0.0000	0.0108	0.0298	0.0000	0.1617	0.7862
	a <sub>12</sub>	1	0.1652	0.1306	0.1566	0.1324	0.1277	0.2874	0.0000
		2	0.0374	0.0923	0.1199	0.4178	0.2719	0.0245	0.0361
		3	0.0467	0.0112	0.0194	0.0279	0.0891	0.2298	0.5759
	a <sub>13</sub>	1	0.0000	0.0000	0.0340	0.1150	0.4882	0.3015	0.0612
		2	0.0593	0.0603	0.0840	0.2334	0.1514	0.2136	0.1981
		3	0.0000	0.0066	0.0068	0.0080	0.0420	0.1263	0.8102
	a <sub>14</sub>	1	0.0000	0.0873	0.0000	0.0000	0.1181	0.7481	0.0464
		2	0.0000	0.0000	0.1482	0.2964	0.2055	0.0922	0.2577
		3	0.0137	0.0000	0.0000	0.0000	0.0206	0.0493	0.9164
	a <sub>15</sub>	1	0.0293	0.0582	0.0287	0.000	0.2635	0.5319	0.0884
		2	0.0000	0.1183	0.1190	0.321	0.2534	0.1552	0.0331
		3	0.0411	0.0069	0.0069	0.008	0.0301	0.1336	0.7734
	a <sub>16</sub>	1	0.0873	0.0291	0.1491	0.0280	0.5203	0.1862	0.0000
		2	0.0852	0.0896	0.1175	0.4252	0.1530	0.0991	0.0305
		3	0.0831	0.0272	0.0268	0.0870	0.1574	0.1663	0.4521
	a <sub>17</sub>	1	0.0000	0.0574	0.0291	0.0000	0.1792	0.4605	0.2738
		2	0.0000	0.0342	0.0891	0.3966	0.1773	0.1848	0.1179
		3	0.0274	0.0540	0.0137	0.0454	0.0539	0.1024	0.7033
	a <sub>18</sub>	1	0.0000	0.0584	0.0000	0.0300	0.2379	0.6445	0.0292
		2	0.1499	0.0904	0.1196	0.3279	0.1023	0.1529	0.0570
		3	0.0544	0.0476	0.0683	0.0199	0.1259	0.1967	0.4871
	a <sub>19</sub>	1	0.0000	0.0291	0.0000	0.0000	0.0629	0.5629	0.3450
		2	0.0296	0.0000	0.0889	0.1484	0.3171	0.2772	0.1387
		3	0.0000	0.0000	0.0000	0.0342	0.0147	0.1872	0.7639
X <sub>5</sub>	a <sub>20</sub>	1	0.0341	0.1000	0.1601	0.0343	0.2779	0.3065	0.0872
		2	0.0709	0.0000	0.0000	0.4495	0.3483	0.1313	0.0000
		3	0.0026	0.0075	0.0000	0.0000	0.0688	0.1309	0.7902
	a <sub>21</sub>	1	0.0732	0.1212	0.1413	0.1166	0.2538	0.2486	0.0453
		2	0.1459	0.0478	0.0821	0.3367	0.2257	0.1367	0.0252
		3	0.0100	0.0000	0.0000	0.0357	0.0662	0.2535	0.6346
	a <sub>22</sub>	1	0.0000	0.1601	0.0250	0.0000	0.2837	0.2869	0.2443
		2	0.0585	0.0000	0.0000	0.4485	0.0560	0.2388	0.1982
		3	0.0201	0.0000	0.0441	0.0218	0.0440	0.1682	0.7018
	a <sub>23</sub>	1	0.2360	0.1806	0.1702	0.1703	0.2429	0.0000	0.0000
		2	0.2494	0.0636	0.0000	0.4142	0.1439	0.1289	0.0000
		3	0.1249	0.0843	0.0931	0.1718	0.2534	0.1038	0.1688
X <sub>6</sub>	a <sub>24</sub>	1	0.0738	0.0927	0.2772	0.1631	0.3185	0.0321	0.0426
		2	0.1474	0.1395	0.0000	0.0000	0.0000	0.2340	0.4791
		3	0.0000	0.0166	0.0000	0.2432	0.3924	0.3096	0.0382
	a <sub>25</sub>	1	0.0000	0.1601	0.0774	0.3035	0.1927	0.1962	0.0701
		2	0.2021	0.0269	0.0198	0.0000	0.1330	0.0943	0.5238
		3	0.0000	0.0219	0.0827	0.1665	0.2039	0.2973	0.2277

Indic- tors	Latent classes	Categories						
		1	2	3	4	5	6	7
$a_{27}$	1	0.0000	0.1997	0.1399	0.4699	0.1613	0.0000	0.0292
	2	0.1604	0.0582	0.1093	0.0699	0.0591	0.2006	0.3425
	3	0.0612	0.0000	0.0796	0.0207	0.3813	0.3772	0.0799
$a_{29}$	1	0.7021	0.0372	0.0000	0.0085	0.0293	0.0289	0.1941
	2	0.1467	0.1417	0.0977	0.2507	0.1579	0.1515	0.0538
$X_7$	$a_{30}$	1	0.4618	0.0000	0.0236	0.0000	0.0074	0.507
		2	0.0590	0.1173	0.0950	0.2802	0.2121	0.202

Table 4 provides the estimated conditional item response probabilities for each of the indicator variables corresponding to the different latent variables. The rows of Table 4 correspond to different latent classes of each latent variables and columns correspond to different categories of each of the indicator variable. For latent variable  $X_1$  (work efficiency), the conditional probabilities  $P[a_1 = 7|X_1 = 1] = 0.6466$  and  $P[a_4 = 7|X_1 = 1] = 0.6393$ , by considering grouping variable gender, envisaged that the respondents (students) are confused on the importance of cell phones for efficient time usage and multitasking given their work efficiency, respectively. While the fact that the same conditional probabilities by considering two grouping variables, namely gender and brand, reduces to 0.0413 and 0.2683, respectively, indicates that the brand of mobile phones plays an important role for students to make their work efficiently.

For latent variable  $X_2$  (up to date), the conditional probabilities obtained by considering grouping variables, gender and brand, did not depict any peculiarities in the responses of the respondents. This means that the respondents believe that they remain up to date with the help of cell phones.

For latent variable  $X_3$  (superiority over landlines), the conditional probabilities  $P[a_9 = 7|X_3 = 1] = 0.8282$ ,  $P[a_{10} = 7|X_3 = 1] = 0.7734$ ,  $P[a_{11} = 7|X_3 = 1] = 0.7881$  and  $P[a_{12} = 7|X_3 = 1] = 0.6037$ , without considering any covariates, indicate that the respondents are not sure about their preference of cell phones over landlines. While the inclusion of the grouping variable gender, these conditional probabilities reduce to 0.1396, 0.0921, 0.0682 and 0.000, respectively, which indicates that the respondents have influence on their preference of cell phones over landlines.

For latent variable  $X_4$  (safety/security), none of the conditional probabilities, when considered with the grouping variable, gender, indicates any sort of unexpected behaviour. This implies that respondents believed that carrying cell phones make them feel safe and secure.

Further, for latent variables  $X_5$  (dependency) and  $X_6$  (negatives), the respondents believe that they are dependent on the cell phones which have negative impacts on their life. And for latent variable  $X_7$  (functionality), the results signify that the respondents are comfortable for providing their opinion about the non-calling functionality of cell phones.

**Table 5:** Estimated coefficients on the covariates along with standard error

Latent variables	Latent classes	Covariates	Coefficient	Std. error	t-statistics	p-value
X <sub>1</sub>	2	(Intercept)	-6.2089	2.2710	-2.734	0.007
		g	4.2355	1.7516	2.418	0.017
		b	2.3226	0.8163	2.845	0.005
		g:b	-1.6652	0.6313	-2.638	0.009
	3	(Intercept)	30.4070	0.4634	65.611	0.000
		g	-31.2964	0.5003	-62.555	0.000
		b	-10.0995	0.3870	-26.097	0.000
		g:b	10.1252	0.2639	38.367	0.000
X <sub>2</sub>	2	(Intercept)	51.2779	4.3713	11.730	0
		g	-43.2907	4.4161	-9.803	0
		b	-18.1251	1.6642	-10.891	0
		g:b	14.7831	1.4530	10.174	0
	3	(Intercept)	46.5247	4.3529	10.688	0.000
		g	-39.9030	4.3566	-9.159	0.000
		b	-15.4059	1.4736	-10.454	0.000
		g:b	13.2461	1.4492	9.140	0.000
X <sub>3</sub>	2	(Intercept)	-1.1287	0.7469	-1.511	0.133
		g	0.6227	0.5670	1.098	0.274
	3	(Intercept)	-0.7873	0.6032	-1.305	0.194
		g	0.8256	0.4564	1.809	0.073
X <sub>4</sub>	2	(Intercept)	-15.3427	0.8012	-19.148	0
		g	14.9957	0.4724	31.743	0
	3	(Intercept)	-14.3230	0.6323	-22.65	0.000
		g	15.3336	0.5957	25.74	0.000
X <sub>5</sub>	2	(Intercept)	-37.4900	0.5915	-63.379	0
		g	35.0006	0.6209	56.369	0
		a	12.8783	0.4421	29.126	0
		g:a	-12.0000	0.3340	-35.919	0
	3	(Intercept)	-5.8590	0.6810	-8.603	0.000
		g	1.6838	0.6976	2.414	0.017
		a	2.0573	0.3646	5.642	0.000
		g:a	-0.4166	0.2767	-1.505	0.135
X <sub>6</sub>	2	(Intercept)	30.3755	0.4583	66.277	0
		g	-31.9703	0.4712	-67.842	0
		a	-10.3418	0.3147	-32.857	0
		g:a	10.9406	0.2395	45.670	0
	3	(Intercept)	29.2397	0.5782	50.568	0.000
		g	-31.2611	0.5784	-54.044	0.000
		a	-9.9287	0.3416	-29.063	0.000
		g:a	10.7219	0.2797	38.331	0.000

Latent variables	Latent classes	Covariates	Coefficient	Std. error	t-statistics	p-value
$X_7$	2	(Intercept)	90.1094	0.2250	400.318	0.000
		g	-42.1254	0.1856	-226.881	0.000
		b	-40.2853	0.1970	-204.410	0.000
		a	-30.1543	0.6356	-47.442	0.000
		g:b	24.9592	0.2027	123.114	0.000
		g:a	14.1658	0.5325	26.598	0.000
		b:a	13.4523	0.2485	54.115	0.000
		g:b:a	-8.2771	0.2005	-41.277	0.000

Table 5 provides the estimated coefficients along with standard errors. It is observed from Table 5 that all the covariates for the selected latent class model with 3 latent classes are highly significant at 5% level of significance. Here, the 1<sup>st</sup> latent classes have been identified as "improve", 2<sup>nd</sup> latent class has been named as "same" and finally the 3<sup>rd</sup> latent class has been named as "worsen", for all latent variables, except for functionality for which the two classes are "non-calling" and "calling", respectively. It is also observed that with inclusion of significant covariates for each of the latent variable, the value of the standard errors and t- statistics significantly reduces. Next, we consider the latent regression model which permits the inclusion of covariates to predict individual latent class membership. The following table provides the log-ratio prior probability that a respondent will belong to the same group with respect to the 1<sup>st</sup> group.

From appendix B, we can calculate the predicted prior probabilities on substituting the response of the  $i^{th}$  individual to the corresponding grouping variable. These probabilities will help to evaluate the estimated latent class membership of individuals with the inclusion of covariates.

## 5. Conclusions

This study considered the problem of selection of indicator variables for different latent variables and to identify different behavioural classes among latent variables. We adopt BGSA (Backward Greedy Search Algorithm) for indicator variables selection and LCA (Latent Class Analysis) for identifying the different behavioural classes. As an application we have investigated the different behaviours of youth towards the usage of mobile phones through questionnaire comprised of 31 items. The questionnaire consists of qualitative questions on usage of cell phones, necessity, cost factors over landlines, safety reasons for carrying cell phones, etc.

The different behaviours of mobile usage has been identified through BGSA as work efficiency, up to date, Superiority over landlines, Safety/Security, Dependency, Negatives and Functionality. Further, we have performed LCA to examine the heterogeneity among the population of youths. The latent classes have been identified as "improve", "same" and "worsen", for all latent variables, except for functionality, which includes two classes, namely "non-calling" and "calling", respectively.

It is envisaged that the different brands availability decreases their efficiency, 45% of students believe that they get updated with cell phones. Also, around 68% of the students

feel safe, while having cell phones and 47% of the students get addicted towards the usage of cell phones. Lastly, due to the non-calling functions of cell phones, 72% of the students gets addicted to cell phones because of the brand of cell phones.

Next, we constructed latent class regression models for each of the latent variable with covariates, which help us to predict the latent class membership of an individual. As the study is limited to the number of respondents with specific specialization, we recommend the proposed methodology to be used for a more diversified sample.

## Acknowledgements

The authors are highly grateful to the learned referees for their constructive comments/suggestions, which helped in the improvement of the paper.

## References

- AKAIKE, H., (1973). Information Theory and an Extension of the Maximum Likelihood Principle. 2<sup>nd</sup> *International Symposium on Information Theory*, pp. 267–281.
- BARTHOLOMEW, D., KNOTT, M. and MOUSTAKI, I., (2011). Latent Variable Models and Factor Analysis. *Wiley*.
- BAUMGARTNER, H. and JAN-BENEDICT, E. M. STEENKAMP, (2006). Response Biases in Marketing Research. *Handbook of Marketing Research*, Thousand Oaks, CA: Sage, pp. 95–109.
- BIEMER, P., (2010). Latent Class Analysis of survey error. *A John Willey and Sons, Inc.* publications.
- BIEMER, P., WIESEN, C., (2002). Latent class analysis of embedded repeated measurements: An application to the National Household Survey on Drug Abuse. *Journal of the Royal Statistical Society, Series A*, 165(1), pp. 97–119.
- BODUSZEK, D., O'SHEA, C., DHINGRA, K. and HYLAND, P., (2014). Latent Class Analysis of Criminal Social Identity in a Prison Sample. *Polish Psychological Bulletin*, 45(2), pp. 192–199.
- DAYTON, C., MACREADY, G. , (1988). Concomitant-Variable Latent-Class Models. *Journal of the American Statistical Association*, 83(401), pp. 173–178.
- DEAN, N., RAFTERY, A. E., (2010). Latent class analysis variable selection. *Annals of the Institute of Statistical Mathematics*, 62, pp. 11–35.

- DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B., (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B (Methodological)*, 39(1), pp. 1–38.
- FOP, M., SMART, K. M. and MURPHY, T. B., (2017). Variable selection for latent class analysis with application to low back pain diagnosis. *Annals of Applied Statistics*, 11(4), pp. 2085–2115.
- FOP, M., MURPHY, T. B., (2017). LCAvarsel: Variable selection for latent class analysis R package version, <https://cran.r-project.org/package=LCAvarsel>.
- FORMANN, A. K., (1984). Constrained latent class models: theory and applications. *British Journal Mathematical and Statistical Psychology*, 38, pp. 87–111.
- FORMANN, A. K., (1992). Linear logistic latent class analysis for polytomous data. *Journal of American Statistical Association*, 87, pp. 476–486.
- GOODMAN, L. A., (1974). Exploratory latent structure analysis using both identifiable and unidentifiable models. *Biometrika*, 61, pp. 215–231.
- HABERMAN, S. J., (1979). Analysis of Qualitative Data. *New York: Academic Press 1979*; Vol. 2: New Developments.
- HAGENAARS, J. A. and MCCUTCHEON, A. L., (2002). Applied Latent Class Analysis. *Cambridge University Press*, New York.
- HUI, S. L., WALTER, S. D., ( 1980 ). Estimating the error rates of diagnostic tests. *Biometrics*, 36, pp. 167–171.
- KASS, R. E., RAFTERY, A. E., (1995). Bayes factors. *Journal of the American Statistical Association*, 90, pp. 773–795.
- KUMAR, S., (2015). Diagnose response bias and heterogeneity: A Latent class approach on Indian household inflation expectation survey. *International Journal of Advances in Social Sciences*, 3(4), pp. 152–158.
- KUMAR, S., (2016). Latent class analysis for reliable measure of inflation expectation in the Indian public. *European Journal of Economic and Statistics*, 1(1), pp. 9–16.
- KUMAR, S., HUSAIN, Z., and MUKHERJEE, D., (2017). Assessing consistency of consumer confidence data using latent class analysis with time factor. *Economic Analysis & Policy*, 55, pp. 35–46.

- LANZA, S. T., RHOADES, B. L., (2013). Latent class analysis: an alternative perspective on subgroup analysis in prevention and treatment. *Prevention science : The Official Journal of the Society for Prevention Research*, 14(2), pp. 157–168.
- LAZARSFELD, P. F., (1950). The logical and mathematical foundation of latent structure analysis and the interpretation and mathematical foundation of latent structure analysis. S.A. Stouffer et al. (eds.), *Measurement and Prediction*, pp. 362–472. Princeton, NJ: Princeton University Press.
- LINZER, D. A., LEWIS, J., (2011). poLCA: Polytomous Variable Latent Class Analysis. *Annals of Applied Statistics*, 11(4), pp. 2085–2115, <http://CRAN.R-project.org/package=poLCA>.
- MCLACHLAN, G. and PEEL, D., (2000). Finite Mixture Models. *John Wiley & Sons*, New York.
- MCLACHLAN, G., KRISHNAN, T., (2008). The EM Algorithm and Extensions. *Wiley*.
- MOOIJART, A. B., (1992). The EM algorithm for latent class analysis with equality constraints. *Psychometrika*, 57(2), pp. 261–269.
- PETERSEN, K. J., QUALTER, P. and HUMPHREY, N., (2019). The Application of Latent Class Analysis for Investigating Population Child Mental Health: A Systematic Review. *Frontiers in Psychology*, Vol.10.
- PORCU, M., GIAMBONA, F., (2017). Introduction to Latent Class Analysis with Applications. *The Journal of Early Adolescence*, 37(1), pp. 129–158.
- RAFTERY, A. E., DEAN, N., (2006). Variable selection for model-based clustering. *The Journal of American Statistical Association*, 101, pp. 168–178.
- SAPOUNIDIS, T., STAMOVLASIS, D. and DEMETRIADIS, S., (2019). Latent Class Modeling of Children Preference Profiles on Tangible and Graphical Robot Programming. *IEEE Transactions on Education*, 62(2), pp. 127–133.
- TRAI, (2018). <https://www.medianama.com/2018/03/223-india-had-1-012-billion-active-mobile-connections-in-january-2018-trai/>.
- VERMUNT, K., JEROEN, K., (2010). Latent class modelling with covariates: Two improved three-step approaches. *Political Analysis*, 18, pp. 450–469.



## APPENDICES

### 6. Detailed description of variables

Variables	Descriptions
a1	A cell phone allows me to use my time efficiently
a2	I use my cell phone to make use of time that otherwise would be wasted
a3	We need a cell phone to be successful in the world today
a4	A cell phone allows me to do two things at once
a5	Those people who do not have a cell phone are out of touch with modern world
a7	I often use my cell phone to schedule/reschedule an appointment at the last minute
a8	It is financially beneficial to use a cell phone as opposed to a landline
a9	A cell phone is more affordable than a landline phone service
a10	If I had to choose, I would use a cell phone instead of a landline because a cell phone is cheaper
a11	A cell phone is a cheaper alternative for long distance calls than a landline
a12	I do not use landlines because having a cell phone is cheaper
a13	Having a cell phone makes me feel safe while I am walking alone at night
a14	My parent wanted me to have a cell phone so I can get in touch with her/him if necessary
a15	I use my cell phone to keep my parent from worrying about me
a16	Having a cell phone makes me feel safe while I am driving
a17	I got my cell phone to use in case of emergency
a18	My parent worries about me less because I have a cell phone
a19	With a cell phone I can keep in touch with my family members
a20	When I do not have my cell phone with me, I feel disconnected
a21	I feel lost when I leave my cell phone at home
a22	I always leave my cell phone on
a23	I feel upset when I miss a call to my cell phone
a24	A cell phone distracts me from being aware of my surroundings
a25	I feel embarrassed by my cell phone ringing at inappropriate times
a26	I am often distracted by my cell phone when driving
a27	I am tired of being accessible all the time
a29	I do not care to learn how to use non-calling functions on my cell phone
a30	I seldom use non-calling functions of my cell phone
g	Gender
b	The brand of a cell phone is important to me
a	A cell phone is addictive

7. Log ratio and Predicted prior probabilities for each of the latent variable

r	Log ratio prior probability $\ln(\frac{p_{ri}}{p_{li}})$	Predicted prior probability $p_{ri}$
$X_1$	2	$\frac{\exp^{-6.2089+4.23554\times g_i+2.3226\times b_i-1.6652\times (g_i:b_i)}}{1+\exp^{-6.2089+4.23554\times g_i+2.3226\times b_i-1.6652\times (g_i:b_i)}+\exp^{30.4070-31.2964\times g_i-10.0995\times b_i+10.1252\times (g_i:b_i)}}$
	3	$\frac{\exp^{30.4070-31.2964\times g_i-10.0995\times b_i+10.1252\times (g_i:b_i)}}{1+\exp^{-6.2089+4.2355\times g_i+2.3226\times b_i-1.6652}+\exp^{30.4070-31.2964\times g_i-10.0995\times b_i+10.1252\times (g_i:b_i)}}$
	2	$\frac{\exp^{51.2779-43.2907\times g_i-18.1251\times b_i+14.7831\times (g_i:b_i)}}{1+\exp^{51.2779-43.2907\times g_i-18.1251\times b_i+14.7831\times (g_i:b_i)}+\exp^{46.5247-39.9030\times g_i-15.4059\times b_i+13.2461\times (g_i:b_i)}}$
$X_2$	3	$\frac{\exp^{46.5247-39.9030\times g_i-15.4059\times b_i+13.2461\times (g_i:b_i)}}{1+\exp^{51.2779-43.2907\times g_i-18.1251\times b_i+14.7831\times (g_i:b_i)}+\exp^{46.5247-39.9030\times g_i-15.4059\times b_i+13.2461\times (g_i:b_i)}}$
	2	$\frac{\exp^{-1.1287+0.6227\times g_i}}{1+\exp^{-1.1287+0.6227\times g_i}+\exp^{-0.7873+0.8256\times g_i}}$
	3	$\frac{\exp^{-0.7873+0.8256\times g_i}}{1+\exp^{-1.1287+0.6227\times g_i}+\exp^{-0.7873+0.8256\times g_i}}$
$X_3$	2	$\frac{\exp^{-15.3427+14.9957\times g_i}}{1+\exp^{-15.3427+14.9957\times g_i}+\exp^{-14.3230+15.3336\times g_i}}$
	3	$\frac{\exp^{-14.3230+15.3336\times g_i}}{1+\exp^{-15.3427+14.9957\times g_i}+\exp^{-14.3230+15.3336\times g_i}}$
	2	$\frac{\exp^{-14.3230+15.3336\times g_i}}{1+\exp^{-15.3427+14.9957\times g_i}+\exp^{-14.3230+15.3336\times g_i}}$

<b>r</b>	<b>Log ratio prior probability <math>\ln(\frac{p_{ri}}{p_i})</math></b>	<b>Predicted prior probability <math>p_{ri}</math></b>
X <sub>5</sub>	2	$\frac{\exp^{-37.4900 + 35.0006 \times g_i + 12.8783 \times a_i - 12.0000 \times (g_i : a_i)}}{1 + \exp^{-37.4900 + 35.0006 \times g_i + 12.8783 \times a_i - 12.0000 \times (g_i : a_i)} + \exp^{-5.8590 + 1.6838 \times g_i + 2.0573 \times a_i - 0.4166 \times (g_i : a_i)}}$
	3	$\frac{\exp^{-5.8590 + 1.6838 \times g_i + 2.0573 \times a_i - 0.4166 \times (g_i : a_i)}}{1 + \exp^{-37.4900 + 35.0006 \times g_i + 12.8783 \times a_i - 12.0000 \times (g_i : a_i)} + \exp^{-5.8590 + 1.6838 \times g_i + 2.0573 \times a_i - 0.4166 \times (g_i : a_i)}}$
X <sub>6</sub>	2	$\frac{\exp^{30.3755 - 31.9703 \times g_i - 10.3418 \times a_i + 10.9406 \times (g_i : a_i)}}{1 + \exp^{30.3755 - 31.9703 \times g_i - 10.3418 \times a_i + 10.9406 \times (g_i : a_i)} + \exp^{29.2397 - 31.2611 \times g_i - 9.9287 \times a_i + 10.7219 \times (g_i : a_i)}}$
	3	$\frac{\exp^{29.2397 - 31.2611 \times g_i - 9.9287 \times a_i + 10.7219 \times (g_i : a_i)}}{1 + \exp^{30.3755 - 31.9703 \times g_i - 10.3418 \times a_i + 10.9406 \times (g_i : a_i)} + \exp^{29.2397 - 31.2611 \times g_i - 9.9287 \times a_i + 10.7219 \times (g_i : a_i)}}$
X <sub>7</sub>	2	$\frac{\exp^{90.1094 - 42.1254 \times g_i - 40.2853 \times b_i - 30.1543 \times a_i} \times (g_i : b_i) + 14.1658 \times (g_i : a_i) + 13.4523 \times (b_i : a_i) - 8.2771 \times (g_i : b_i : a_i)}}{1 + \exp^{90.1094 - 42.1254 \times g_i - 40.2853 \times b_i - 30.1543 \times a_i} + 24.9592 \times (g_i : b_i) + 14.1658 \times (g_i : a_i) + 13.4523 \times (b_i : a_i) - 8.2771 \times (g_i : b_i : a_i)}}$
	3	$\frac{\exp^{90.1094 - 42.1254 \times g_i - 40.2853 \times b_i - 30.1543 \times a_i} \times (g_i : b_i) + 14.1658 \times (g_i : a_i) + 13.4523 \times (b_i : a_i) - 8.2771 \times (g_i : b_i : a_i)}}{1 + \exp^{90.1094 - 42.1254 \times g_i - 40.2853 \times b_i - 30.1543 \times a_i} + 24.9592 \times (g_i : b_i) + 14.1658 \times (g_i : a_i) + 13.4523 \times (b_i : a_i) - 8.2771 \times (g_i : b_i : a_i)}}$

where r: latent class (r=2,3) and

i: no. of respondents (i=1,2,...,N)