# 國立交通大學

## 統計學研究所

## 碩 士 論 文

多維營養資料之因子
及長期追蹤分析

Factor and Longitudinal Analysis for

Multiple Nutritional Data

研 究 生：劉育倫

指導教授：黃冠華 博士

中 華 民 國 九 十 五 年 六 月

多維營養資料之因子
及長期追蹤分析

# Factor and Longitudinal Analysis for Multiple Nutritional Data

研 究 生：劉育倫　　　　　Student : Yu-Lun Liu

指導教授：黃冠華　　　　　Advisor : Dr. Guan-Hua Huang

國 立 交 通 大 學

統計學研究所

碩 士 論 文

A Thesis

Submitted to Institute of Statistics

College of Science

Nation Chiao-Tung University

in partial Fulfillment of the Requirements

for the degree of Master

in

Statistics

June 2006

Hsinchu, Taiwan

中華民國九十五年六月

# 多維營養資料之因子及長期追蹤分析

學生：劉育倫　　　　指導教授：黃冠華 博士

## 國立交通大學統計學研究所

## 摘　要

此論文主要分兩個部份去研究營養的相關問題。第一部份,主要是探討美國社區老人的口腔健康和營養間的關係。營養不良普遍發生在老年人中,長期營養不良可能會引發疾病而造成身體不適,近而影響生活的品質。然而老年人常因為牙齒的因素而無法進食而導致營養不良,因此口腔健康是不容忽視的。利用因子分析建立口腔健康和營養之間可能存在的關係。第二部份,針對台大醫院住院的老人做營養的長期追蹤,分別在出院前和出院後三到六個月做調查。主要是觀察住院老人的營養如何隨著時間的變化而有所改變, 應用因子分析和廣義估計方程式去分析營養變化,且找出影響其變化的因素如個人生理或者心理社會因子等。

關鍵字：因子分析 ；廣義估計方程式

# Factor and Longitudinal Analysis for Multiple Nutritional Data

Student : Yu-Lun Liu            Advisor : Dr. Guan-Hua Huang

Institute of Statistics

National Chiao Tung University

## Abstract

This thesis consists of two parts which are nutrition related problems. In the first part, we study the relationships between oral health and nutrition for community-dwelling in U.S.A.. Malnutrition is common in elder people, it may influence the elder people' health and the life quality. However, elder people are usually unable to eat due to teeth, it results in poor nutrition. Thus, oral health is also important component for elder people's health. We use factor analysis to obtain the possible associations between oral health and nutrition. In the second part, it is a longitudinal data set for hospitalized elder people of National Taiwan University Hospital, we collect the data in different time points as follows: before discharge, three month post index hospitalization and six month post index hospitalization. Factor analysis and Generalized estimating equations (GEEs) are employed to study the change of the nutritional status over time and find the risk factors such as personal physical factors, psychosocial factors and so on.

*Key words*: Factor analysis; Generalized estimating equations

# 誌謝

撰寫論文的過程中，因為有許多人的幫忙才能順利地完成。在此誠心地感謝大家，倘若沒有大家的協助可能無法完成此論文。首先，最要感謝我的指導教授-黃冠華老師，他很有耐心地指導每一個步驟。當我遇到瓶頸時，他會分析問題的癥結，引導我如何去解決問題及提供意見，使我受用無窮。還有，陳佳慧老師所提供的資料及一些問題的探討，使論文有發揮的空間。此外，還要感謝口試委員陳鄰安老師和許文郁老師，他們提供一些建議及想法，並匡正論文錯誤的地方,使得它能夠更完整。

另外，無論精神上的勉勵或是實質上的協助，由衷地感謝我的同學和朋友以及祝福大家。最後，感謝我的家人和親戚，一路走來他們的鼓勵及支持是督促我不斷努力的動力，讓我更加有信心達成目標，並謹以此成果獻給他們。
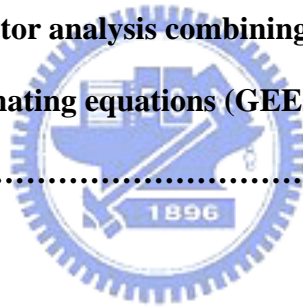
劉育倫　　謹誌于

國立交通大學統計學研究所

中華民國九十五年六月

# Contents

# List of Tables

# List of Figures

# Chapter 1

# INTRODUCTION

Malnutrition is prevalent in elder people, so it is not overlooked for elder people's general health. In this article, we are interested in discovering the nutritional status related problems, thus there are two main questions. First, what is the relationship between nutrition and oral health? And second, what are the nutritional changes of hospitalized elderly patients over time.

Nutrition and oral health are measured by multiple indicators. Previous study uses the total sums of nutrition and oral health indicators to study their relationship. This may mask the association. Here, we use factor analysis to find underlying structure of nutrition and oral health, thus they can be more accurate. The former papers have many discussions on oral health and nutrition respectively, but here we put two together to find the possible relationships. In this thesis, we are interested in discovering the relationships of oral health and nutrition for community-dwelling elders

and we also add some risk factors including comorbidity, depressive symptoms, satisfaction with support and so on. By adding risk factors, we can construct the relationships between oral health and nutrition excluding the effects due to confounders.

Nutrition is measured by multiple indicators. Previous study uses the sum of these indicators to represent the nutrition status. This might mass the true nutrition. We use factor analysis to draw the distinct dimension of nutrition. Previous studies have not collected the longitudinal data of nutrition status. The main advantage of a longitudinal study is its effectiveness for studying change. Another merit of the longitudinal study is its ability to distinguish the degree of variation in response variables across time for one person from the variation in response variables among people (Digglen *et al.*,2002).Thus, longitudinal data enable us to study the association between nutrition and risk factors at different stage of hospitalization, and to study the nutrition change and its associated causes.

In the second part of this thesis, we combine factor analysis and longitudinal data analysis to study the nutritional change for different dimension of nutrition. Using longitudinal data, our study examined factors that are associated with nutritional change over time in hospitalized elderly persons and investigated interactions among those risk factors. In this work we applied generalized estimating equations (GEEs) for the analysis of nutritional change for hospitalized older people. The advantage of the GEE approach is that it requires weaker distributional assumptions and maintains the properties of consistency and asymptotic normality of parameter

estimates. And useful of the GEE approach is that it is not necessary for the "working" correlation matrix to be correctly specified to construct consistency and asymptotic normality of parameter estimates (Zeger and Liang,1986).

# Chapter 2

# LITERATURE REVIEW

This chapter reviews some useful statistical tools for modeling data that are used by this thesis and are implemented in Mplus Version3.13 and R package. First, latent variable models and factor analysis with categorical indicators are illustrated with community-dwelling elder people's health studies. The factor analysis has discussed by several authors (Reyment and Joreskog ,1993; Basilevsky,1994). It is generally understood to refer to a set of closely related models intended for exploring and establishing correlation structure among the observed variables. In short words, the factor analysis uses the fact that measured variables can be correlated in such a way that their correlation may be reconstructed by a smaller set of factors, which could represent the underlying structure in a concise and interpretable form.

In Mplus ( Muthén and Muthén,1998-2005), there are two types of factors analysis: exploratory factor analysis (EFA) and confirmatory factor

analysis (CFA). EFA is used to determine the number of continuous latent variables (or factors) that are needed to explain the correlations among a set of observed variables (or factor indicators). The goal of EFA is to find the smallest number of interpretable factors that can adequately explain the correlations among a set of variables. Furthermore, CFA is used to describe the relationships between a set of observed variables and a set of continuous latent variables. It enables the investigators to theorize an underlying structure and justify whether the observed data fit this a prior hypothesized model. Traditional EFA and CFA are for continuous observed variables. For categorical observed variables, the liability threshold model is used, which postulates the existence of an unobserved continuous variable and a set of thresholds.

In Mplus, it also provided several fit indices to assess the performance of data-model fits. CFA are used frequently in preparation for analyzing more general structural equation models (SEM) on the side. CFA plays an important role in structural equation modeling. Structural equation models (SEMs), also called simultaneous equation models, are multivariate (i.e., multi-equation) regression models. Unlike the more traditional multivariate linear model, however, the response variable in one regression equation in an SEM may appear as a predictor in another equation; indeed, variables in an SEM may influence one-another reciprocally, either directly or through other variables as intermediaries. These structural equations are meant to represent causal relationships among the variables in the model. In fact, SEM is composed of two parts: the measurement model and the structural

model. The measurement model relates observed indicators to latent variables and sometimes to observed covariates. The structural model then specifies relations among latent variables and regressions of latent variables on observed variables. When the observed indicators are categorical, we need to modify the conventional measurement model for continuous indicators. However, the structural model can remain essentially the same as in the continuous case (Skrondal,2005). SEM may be used as a more powerful alternative to multiple regression, path analysis, factor analysis, time series analysis, and analysis of covariance. That is, these procedures may be seen as special cases of SEM, or, to put it another way, SEM is an extension of the general linear model (GLM) of which multiple regression is a part. Bollen (1989) provide a introduction to the general structural equation system and emphasize the application of techniques. Also, SEM has emerged as a helpful multivariate data analysis tool in social science research settings, especially in the fields of sociology, psychology, and education (Mueller ,1996).

Secondly, many researchers are interested in analyzing data from longitudinal studies whose feature is that individuals are measured repeatedly through time. Zeger and Liang (1992) discuss the statistical methods for the analysis of discrete and continuous longitudinal data using three approaches, marginal, transition and random effect models. Liang and Zeger (1986) describe the marginal expectation of the outcome variable as a function of the covariates while accounting for the correlation among the repeated observations for a given subject. In addition, they specify a "work-

ing" correlation matrix for the observations for each subject. These authors also formalized an approach to this longitudinal data using generalized estimating equations (GEEs) to extend generalized linear models (GLMs) to a regression setting with correlated observations within subjects. This setup leads to generalized estimating equation (GEEs) which give consistent estimators of the regression coefficients and of their variances under weak assumptions about the actual correlation among a subject's observations. In general, GEEs are used to characterize the marginal expectation of a set of outcomes as a function of a set of study variables. In a marginal model, the analyst is interested in modeling the marginal expectation (average response for observations sharing the same covariates) as a function of explanatory variables. However, in this article we are instead of the expectation of factor scores as dependent variables. Nicholas and Stuart (1999) compare the GEEs implementations of several general purpose statistical packages including SAS, Stata, SUDAAN, and S-Plus. In this paper, we utilize an R package to analyze the nutritional status of hospitalized elders using GEEs approach. The use of GEEs to estimate regression coefficients specified by marginal models has been studied extensively over the last fifteen years. For more detailed treatments, see Prentice (1988), Zhao and Prentice(1990), Thall and Vail(1990), Liang *et al.*(1992), and Fitzmaurice and Laird(1993). In the followings, we will describe the detailed statistical model of factor analysis and GEE approach.

## 2.1 Factor analysis

(1)Variables

Let $\mathbf{y}_i = (y_{i1}, ..., y_{ip})$(dependent variables) be $p$-dimensional categorical variables of individual corresponding to continuous latent response variable $y_i^*$. Let $i(i = 1, 2, ..., n)$ denote the observational unit (the individual) and $j(j = 1, 2, ..., p)$denote the observed dependent variable. A categorical variable $y_{ij}, \ j = 1, 2, ..., p$ with $C$ ordered categories is defined as

$$y_{ij} = c \text{ , if } \tau_{j,c} < y_{ij}^* \leq \tau_{j,c+1} \text{ , } i = 1, 2, ..., n \qquad (2.1)$$

for categories $c = 0, 1, 2, ..., C - 1$and $\tau_0 = -\infty$, $\tau_c = \infty$.

(2)Model

The exploratory factor analysis model is defined as

$$\mathbf{y}_i^* = \mathbf{v} + \mathbf{\Lambda}\boldsymbol{\eta}_i + \boldsymbol{\varepsilon}_i \text{ ,} \qquad (2.2)$$

where $\mathbf{y}^*$ is a $p$-dimensional vector of response variables, $\mathbf{v}$ is a $p$-dimensional parameter vector of measurement intercepts , $\mathbf{\Lambda}$ is a $p \times m$ parameter matrix of factor loadings, $\boldsymbol{\eta}$ is an $m$-dimensional vector of factors, $\boldsymbol{\varepsilon}$ is a $p$-dimensional vector of residual.

The equation (2.2) can be extended as the following to incorporate covariate effects,

$$\mathbf{y}_i^* = \mathbf{v} + \mathbf{\Lambda}\boldsymbol{\eta}_i + \mathbf{K}\mathbf{x}_i + \boldsymbol{\varepsilon}_i \text{ ,} \qquad (2.3)$$

here $\mathbf{K}$ is a $p \times q$ parameter matrix of regression slopes, $\mathbf{x}$ is a $q$-dimensional

vector of independent variables. This equation (2.3) is the measurement part of SEM model, and another part of SEM model is the structural part which is defined in terms of the latent variable regressed on each other and the independent variables,

$$\boldsymbol{\eta}_i = \boldsymbol{\alpha} + \mathbf{B}\boldsymbol{\eta}_i + \boldsymbol{\Gamma}\mathbf{x}_i + \boldsymbol{\zeta}_i \ , \tag{2.4}$$

where $\boldsymbol{\alpha}$ is an $m$-dimensional parameter vector, $\mathbf{B}$ is $m \times m$ parameter matrix of slopes for regressions of latent variables on other latent variables, $\boldsymbol{\Gamma}$ is an $m \times q$ slopes parameter matrix for regressions of the latent variables on the $\mathbf{x}$ variables, and $\boldsymbol{\zeta}$ is an $m$-dimensional vector of residuals(see Muthén1979,1983,1984,1989b).

(3)The scaling parameter of $\boldsymbol{\Delta}$

The scaling parameter $\Delta$ is used as the scale of original response variable $y_i^*$,

$$\mathbf{y}_{si}^* = \boldsymbol{\Delta}\mathbf{y}_i^*.$$

The diagonal elements of $\boldsymbol{\Delta}$ are useful when comparing the same $y$ variables over time, the $\Delta$ element for the first time point can be standardized to one whereas $\Delta$ elements can be estimated for other time points to capture differences in $\mathbf{y}^*$ variances over time.

(4)The threshold parameter of $\boldsymbol{\tau}$

The threshold parameters are used in the model for categorical $y$ variables. As previous equation (2.1), $\tau_{j,c}$ and $\tau_{j,c+1}$ are the threshold pa-

rameters. Given the conditional normality assumption this leads to the
univariate and bivariate probability expressions

$$P(y_i = 1|\mathbf{x}) = \int_{\tau_j^* - \mu_j^*(\mathbf{x})}^{\infty} \phi_1(y_j^*|\mathbf{x}) dy_j^* \ ,$$

$$P(y_j = 1, y_k = 1|\mathbf{x}) = \int_{\tau_j^* - \mu_j^*(\mathbf{x})}^{\infty} \int_{\tau_k^* - \mu_k^*(\mathbf{x})}^{\infty} \phi_2(y_j^*, y_k^*|\mathbf{x}) dy_k^* dy_j^* \ ,$$

where $\tau_j^*$ is the threshold parameter for $y_j^*$ multiplied by the jth diagonal
element of $\mathbf{\Delta}$ , $\phi_1$ is a univariate standard normal density, $\phi_2$ is a bivariate
normal density with unit variances, zero means, and correlation coefficient
$\sigma_{jk}^*$ which is an off-diagonal element of $\mathbf{\Sigma}^*$. The off-diagonal elements of
$\mathbf{\Sigma}^*$ are referred to as probit residual correlations. The elements of $\mathbf{\Delta\Pi}$ are
referred to as probit slopes.

(5)Estimator

For categorical outcomes, the model parameters are estimated by the
mean-adjusted and variance-adjusted weighted least square estimating method
(WLSMV). WLSMV-weighted least square parameter estimates using a di-
agonal weight matrix with robust standard errors and mean-adjusted and
variance-adjusted test statistic. Following is the brief description of this es-
timator. Muthén (1981a) proposed a three-stage limited information WLS
estimator. The three parts are respectively a mean/threshold/reduced-form
regression intercept structure, a reduced-form regression slop structure, and
a covariance/correlation structure.

Consider the three population vectors $\sigma_1$ , $\sigma_2$ and $\sigma_3$ ( Muthén,1983):

10

Part1 contains a mean, or threshold, or intercept structure

$$\sigma_1 = \mathbf{\Delta}^*[\mathbf{K}_\tau \boldsymbol{\tau}_z - \mathbf{K}_v(\mathbf{v}_z + \mathbf{\Lambda}_z(\mathbf{I} - \mathbf{B}_z)^{-1}\boldsymbol{\alpha}_z)] \; ,$$

Part2 contains a slop structure

$$\sigma_2 = vec[\mathbf{\Delta}\mathbf{\Lambda}_z(\mathbf{I} - \mathbf{B}_z)^{-1}\mathbf{\Gamma}_z] \; ,$$

Part3 contains a covariance, correlation, or residual correlation structure

$$\sigma_3 = \mathbf{K}vec\{\mathbf{\Delta}[\mathbf{\Lambda}_z(\mathbf{I} - \mathbf{B}_z)^{-1}\mathbf{\Psi}_z(\mathbf{I} - \mathbf{B}_z)^{'-1}\mathbf{\Lambda}_z' + \mathbf{\Theta}_z]\mathbf{\Delta}\} \; .$$

Here, $\mathbf{\Delta}$ is a diagonal $p \times p$ matrix of scaling factors, $\mathbf{\Delta}^*$contains the same element as $\mathbf{\Delta}$ but diagonal elements are duplicated for categorical variables with more than one threshold, $\mathbf{K}_\tau$ and $\mathbf{K}_v$ similarly distributed elements from the vectors they pre-multiply, the $vec$ operator strings out matrix elements row-wise into a column vector, and $\mathbf{K}$selects lower-triangular elements from the symmetric matrix elements it pre-multiplies, where a diagonal element if only included if the corresponding observed variable is continuous, $\mathbf{\Lambda}_z$ is a $p \times m$ matrix of loadings, $\mathbf{B}_z$ is an $m \times m$ matrix of slopes for the regression among the $m$ latent variables, $\Gamma$ is an $m \times q$ matrix of slopes for the regression of the $m$ latent variables on the $qx$ variables, $\mathbf{\Psi}_z$ is $m \times m$ an covariate matrix for the latent variables and the residuals in the latent variable relations.

With the normality specification on the latent response variables, any model that fits in the general framework is identified if and only if its

parameters are identified in terms of $\boldsymbol{\sigma}^{(1)}$, $\boldsymbol{\sigma}^{(2)}$,...,$\boldsymbol{\sigma}^{(G)}$ , where $\boldsymbol{\sigma}^{(g)'} = (\sigma_1^{(g)'}, \sigma_2^{(g)'}, \sigma_3^{(g)'})$. Muthén (1981a) utilized this fact in that statistics $\mathbf{s}^{(g)}$ were produced as consistent estimators of $\boldsymbol{\sigma}^{(g)}$, in order to estimate the model parameters in a final estimation stage. Preceding estimation stages give $\mathbf{s}^{(g)}$, where only limited information from bivariate distributions is needed. In the final estimation stage , a WLS fitting function with a general, full weight matrix is used

$$F = \sum_{g=1}^{G} (\mathbf{s}^{(g)} - \boldsymbol{\sigma}^{(g)})' \mathbf{W}^{(g)-1} (\mathbf{s}^{(g)} - \boldsymbol{\sigma}^{(g)}) \; ,$$

where the generalized least squares estimator is obtained when $\mathbf{W}^{(g)}$ is a consistent estimator of the asymptotic covariance matrix of $\mathbf{s}^{(g)}$.

(6)Fit indices:

The most commonly used test of model adequacy is the $\chi^2$ goodness-of-test. The null hypothesis for this test is that the model adequately accounts for the data, while the alternative is that there is significant amount of discrepancy. The $\chi^2$ approximation is sensitive to sample size and violation of the multivariate normality assumption. Muthén and Kaplan (1992) proposed that the $\chi^2$ is also sensitive to model complexity. The Mplus provided some model fit measures, the model fit information obtained from these fit indices are very different from that obtained from the $\chi^2$ measure where a hypothesized model is compared to a saturated model. The first fit indices are Tucker-Lewis Index (TLI) and Comparative Fit Index (CFI),

$$TLI = (\chi_B^2/d_B - \chi_{H_0}^2/d_{H_0})/(\chi_B^2/d_B - 1) \ ,$$

$$CFI = 1 - max(\chi_{H_0}^2 - d_{H_0}, 0)/max(\chi_{H_0}^2 - d_{H_0}, \chi_B^2 - d_B, 0) \ ,$$

where $d_B$ and $d_{H_0}$ are the degree of freedom for the baseline and hypothesized models, respectively. TLI, compare the fit of the proposed model to that of a "null model", it has been shown to be much less sensitive to sample size. Hu and Bentler (1999) recommended cutoff values of TLI and CFI both close to 0.95. Both TLI and CFI are incremental fit indices, which measure the improvement of fit by comparing a hypothesized model with a more restricted model. And the second fit index is Root-mean-square Error of approximation (RMSEA). With categorical outcomes, RMSEA is defined as

$$RMSEA = \sqrt{\max[(2F(\widehat{\theta})/d - 1/n), 0]} \ ,$$

where $d$ is a function of the sample variances, and $F(\widehat{\theta})$ is the minimum of the fitting function $F(\theta)$. Browne and Cudeck (1993) suggested that RMSEA values larger than 0.1 are indicative of poor fitting models , values in the range of 0.05 to 0.08 are indicative of fair fit , and values less than 0.05 are indicative of close fit. Also, Hu and Bentler (1999) recommended a cutoff value of RMSEA close to 0.06.

The final fit index is Weighted Root-mean-square residual (WRMR), it is expressed as

$$WRMR = \sqrt{\frac{2nF(\widehat{\theta})}{e}} \ ,$$

where $F(\widehat{\theta})$ is the minimum of the fitting function $F(\theta)$, and $e$ is the number of sample statistics. Small values of WRMR indicate good fit. Yu and Muthén (2001) recommended the value of WRMR $< 0.9$ for good models with categorical outcomes.

## 2.2 Factor scores

Factor scores which are the estimated values of the common factors are not estimates of unknown parameters in the usual sense. Rather, they are estimates of values for the unobserved random factor vectors. Factor scores, an quantify individual cases on a latent continuum using a z-score scale which ranges from approximately -3.0 to +3.0. Bartlett (1937) has suggested that weighted least squares be used to estimate the common factor values.

The weighted least squares method to calculate the factor scores, $\widehat{f}_j$ , as follows

$$\widehat{f}_j = (\widehat{L}'\widehat{\Psi}^{-1}\widehat{L})^{-1}\widehat{L}'\widehat{\Psi}^{-1}(\mathbf{x}_j - \widehat{\mu}) = \widehat{\Delta}^{-1}\widehat{L}'\widehat{\Psi}^{-1}(\mathbf{x}_j - \overline{x}) \ , \quad j = 1, 2, ..., n$$

where $\widehat{\mu}$ denote the estimates of the mean vector, $\widehat{L}$ denote the estimates of the factor loadings, $\widehat{\Psi}$ are the estimates of the specific variances, $\widehat{\Delta} = \widehat{L}'\widehat{\Psi}^{-1}\widehat{L}$ is a diagonal matrix, and $\mathbf{x}_j$ is the observed values.

## 2.3 Generalized Estimating Equations (GEEs)

### 2.3.1 Marginal models

In the regression, we model the marginal expectation, $E(Y_{it})$, as a function of explanatory variables. By marginal expectation, we means the average response over the sub-population that shares a common value of $x$. Assume that

(1)the marginal expectation of the response $Y_{it}$, $E(Y_{it}) = \mu_{it}$ ,depends on explanatory variables, $x_{it}$ , by

$$g(\mu_{it}) = x'_{it}\boldsymbol{\beta} ,$$

where $g$ is a known link function such as the logit for binary responses or log for counts;

(2)the marginal variance depends on the marginal mean according to

$$Var(Y_{it}) = v(\mu_{it})\phi ,$$

where $v$ is a known function and $\phi$ is a scale parameter which may need to be estimated;

(3)the covariance between $Y_{is}$ and $Y_{it}$, $s < t = 1, ..., n_i$ is a function of the marginal means and additional parameters $\alpha$ , that is,

$$cov(Y_{it}, Y_{is}) = c(\mu_{is}, \mu_{it}; \alpha) ,$$

where $c$ is a known function (see Zeger & Liang,1992).

## 2.3.2  Quasi-Likehood

This quasi-likelihood was first proposed by Wedderburn (1974) and later tested extensively by McCullagh (1983). An alternative generalization was proposed by Lee and Nelder (1996, 2001). An extensive review of the development of the GEE approach is given by Ziegler, Kastner, and Blettner (1998).

Assume the observations $(y_{ij}, x_{ij})$ for times $t_{ij}$, $j = 1, 2, ..., n_i$ and subjects $i = 1, 2, ..., K$.

Let $\mathbf{Y}_i = (Y_{i1}, Y_{i2}, ..., Y_{it}, ..., Y_{in_i})'$ and $\mathbf{x}_i$ be $n_i \times p$ matrix $(x_{i1}, x_{i2}, ..., x_{in_i})'$ for the $ith$ subject and the expectations $E(Y_{it}) = \mu_{it}$ be related to the $p$-dimensional regressor $x_{it}$ by the mean-link function $g$

$$g(\mu_{it}) = x_{it}'\boldsymbol{\beta} \ , \tag{2.5}$$

where $\boldsymbol{\beta}$ is a $p \times 1$ vector of parameters. Let

$$Var(Y_{it}) = h_{it}\phi \ , \tag{2.6}$$

where $\phi$ is a scale parameter and $h_{it} = h(\mu_{it})$ is a known variance function. The focus of quasi-likelihood is on methods for inference about $\boldsymbol{\beta}$. Hence, $\phi$ is treated as a nuisance parameter. The quasi-likelihood estimator is the solution of the score-like equation system

$$\mathbf{S}_k(\boldsymbol{\beta}) = \sum_{i=1}^{K} \frac{\partial \boldsymbol{\mu}_i'}{\partial \boldsymbol{\beta}_k} Var_i^{-1}(\mathbf{Y}_i - \boldsymbol{\mu}_i) = 0 \ , \ k = 1, 2, ..., p \tag{2.7}$$

This equation (2.7) are in fact score equation for $\boldsymbol{\beta}$ when $\mathbf{Y}_i$ has distribution from the exponential family. Their solution can be obtained by an iteratively reweighted least squares (see Zeger and Liang,1986). The resulting estimator is asymptotically Gaussian under mild regularity conditions (McCullagh, 1983). It also possesses a Gauss-Markov-like optimally in that is asymptotically the minimum variance estimator among those with linear influence function. Wedderburn (1974) and McCullagh (1983) provide details about quasi-likelihood in the regression context.

### 2.3.3 Generalized Estimating Equations(GEEs)

Zeger and Liang (1986) introduced the GEEs approach for the regression analysis of correlated observations. To utilize the quasi-likelihood approach above, we need to suppose the mean and covariance of the vector of responses,$\mathbf{Y}_i$ , for subject.

Let $\mathbf{R}_i(\boldsymbol{\alpha})$ be the $n_i \times n_i$ working correlation matrix for each $\mathbf{Y}_i$ with the working covariance matrix $\mathbf{V}_i(\boldsymbol{\alpha})$,

$$\mathbf{V}_i(\boldsymbol{\alpha}) = \mathbf{A}_i^{1/2}\mathbf{R}_i(\boldsymbol{\alpha})\mathbf{A}_i^{1/2}/\boldsymbol{\phi} \ ,$$

where $\mathbf{A}_i$ is an $n_i \times n_i$ diagonal matrix with entries $h_{it}$ .

The extension of equation (2.7) is defined by

$$\sum_{i=1}^{K}\mathbf{D}_i^{'}\mathbf{V}_i^{-1}\mathbf{S}_i = \mathbf{0} \ , \tag{2.8}$$

where $\mathbf{S}_i = \mathbf{Y}_i - \boldsymbol{\mu}_i$ with $\boldsymbol{\mu}_i = (\mu_{i1}, ..., \mu_{in_i})^{'}$ and $\mathbf{D}_i = \partial\boldsymbol{\mu}_i/\partial\boldsymbol{\beta}$.

More generally, $\mathbf{U}_i(\boldsymbol{\beta}, \boldsymbol{\alpha}) = \mathbf{D}_i^{'}\mathbf{V}_i^{-1}\mathbf{S}_i$ is equivalent to the estimating function suggested by Wedderburn (1974) except that the $\mathbf{V}_i$'s here are functions of $\boldsymbol{\alpha}$ as well as $\boldsymbol{\beta}$. For any given $\mathbf{R}_i(\boldsymbol{\alpha})$, the estimate, $\widehat{\boldsymbol{\beta}}_R$, of $\boldsymbol{\beta}$ is defined as the solution of

$$\sum_{i=1}^{K} \mathbf{U}_i\{\boldsymbol{\beta}, \widehat{\alpha}[\boldsymbol{\beta}, \widehat{\phi}(\boldsymbol{\beta})]\} = \mathbf{0} \ . \tag{2.9}$$

Under mild regularity conditions, Liang and Zeger (1986) show that $\widehat{\boldsymbol{\beta}}_R$ is a consistent estimator of $\boldsymbol{\beta}$ in the equation (2.9) as $\mathbf{K} \to \infty$. To solve the GEE for $\widehat{\boldsymbol{\beta}}_R$, we iteratively solve for the regression coefficients and the correlation and scale parameters, $\boldsymbol{\alpha}$ and $\boldsymbol{\phi}$. Given an estimate of $\mathbf{R}_i(\boldsymbol{\alpha})$ and of $\boldsymbol{\phi}$, we can calculate an updated estimate of $\boldsymbol{\beta}$ by iteratively reweighted least squares as described by (McCullagh & Nelder,1983).

## 2.3.4 The working correlation

In addition, it is important to specify the working correlation matrix $\mathbf{R}$. There are a variety of common structures including independence, exchangeable, unstructured, auto-regressive, m-dependent, and fixed. Fitzmaurice *et al.*(1993) discuss four common specifications of the working correlation matrix $\mathbf{R}_i(\boldsymbol{\alpha})$ for observations $Y_{is}$ and $Y_{it}$ as follows:

(1)$\mathbf{R}_i(\boldsymbol{\alpha}) = \mathbf{I}$ , where $\mathbf{I}$ is a $n_i \times n_i$ identity matrix. This corresponds to the" working independence" assumption, and gives estimating equations identical to (2.8).

(2) Exchangeable correlation: $corr(Y_{is}, Y_{it}) = \alpha$ for $s \neq t$.

(3) Autoregressive correlation: $corr(Y_{is}, Y_{it}) = \alpha^{|s-t|}$ for $s \neq t$.

(4) Unstructured or pairwise correlation: $corr(Y_{is}, Y_{it}) = \alpha_{st}$ , where $\alpha$ is a $n_i(n_i - 1)/2 \times 1$ vector containing all the pairwise correlations.

Generally speaking, if the number of observations per cluster is small in a balanced and complete design, then an unstructured matrix is recommended (Horton and Lipsitz,1999). Here, we only use the exchangeable working correlation matrix in our analysis.

# Chapter 3

# FACTOR ANALYSIS FOR MULTIPLE NUTRITIONAL DATA

## 3.1 Introduction

Malnutrition is prevalent in elder people. In the United States, it is estimated that 40% of nursing home residents, 50% of hospitalized elders, and 45% of home care elders are malnourished (Nutrition Screening Initiative,1993). And oral health is often neglected component of elder people' health, the tooth loss affects dietary quality and nutrient intake. In general, oral health may influence nutrition, speech, communication and self-image. Furthermore, poor oral health can affect dietary quality and nutrient intake in a manner that potentially increases the risk of several systemic diseases

(Ritchie *et al.*2002). Hence, nutrition and oral health are both important factors for the health of elder people. Because of the two problems, they would affect gradually quality of life for elder people.

Nutrition and oral health status are measured by multiple indicators. Previous study uses the total sums of nutrition and oral health indicators to study their relationship. This may mask the true association. Here, we use factor analysis to find underlying structure of nutrition and oral health, thus they can be more accurate. The former papers have many discussions on oral health and nutrition respectively, but here we put two together to find the possible relationships. In this paper, we are interested in discovering the relationships of oral health and nutrition for community-dwelling elders and we also add some risk factors including comorbidity, depressive symptoms, satisfaction with support and so on. By adding factors, we can construct the relationships between nutrition and oral health excluding the effects due to confounders.

Later, we will present the used statistical methods, including exploratory factor analysis (EFA), confirmatory factor analysis (CFA) and structural equation modeling (SEM) and depict how to collect the sample in section 3.2. Section 3.3 shows the results of the analysis, and in section 3.4 we discussed some findings and provided suggestions for future researches.

## 3.2 Method

### 3.2.1 Sample

This study was designed as a cross-sectional population-based study. Subjects were recruited from an inner-city senior housing complex in Southern Connecticut. The facility was independent living in nature and approximately one-third of the residents were subsidized by public funds. Each resident in the facility was invited to participate. Subjects were excluded if they are unable to give consent (n =5), hospitalized or too ill (n =7), or not English-speaking (n =26). Consent by proxy was not used due to its ethical and methodological difficulties. At the end of a 7-month recruitment period, 243 subjects, from an eligible population of 268 subjects, completed a structured in-home assessment by a trained Geriatric Nurse Practitioner (GNP). In order to alleviate the respondent burden, subjects were giving the option of breaking the interview into two sessions and only two out of 243 subjects used this option. The participation rate was 91%. The non-participants (n =25) were not significantly different from the participants on age ($p = 0.38$), gender($p = 0.23$),marital status ($p = 0.10$),or ethnicity ($p = 0.74$).

### 3.2.2 Risk factors and measured instruments

Chen *et al.*(submitted) proposed possible factors which impacted on malnutrition in the elderly. In that paper, they provided four conceptual texts including loss, chronic illness, dependency and loneliness to test the

22

relationships with the development of malnutrition in the elderly. In loss part, it contains some empirical indices including age, oral health. The chronic part includes as follows: comorbidity, number of medications. For illness, there are several indices including gender/education, functional status. And the loneliness mainly contains social support and depression empirical indices.

Demographics collected in the interview included age, gender, ethnicity, living status, marital status, education and religion.

Oral health was measured by 12-item Geriatric Oral Health Assessment Index (GOHAI). The GOHAI assessed the dimensions of function (eating and speaking), pain, discomfort, worry, and social functioning. In the questionnaire, each item states an oral health related problem and participants are asked to indicate three choices (always/sometimes/never) of how they feel the way described. The sum scores range from 12 to 36 with a high value indicating better self perceived oral health (Atchison and Dolan, 1990).

Nutritional status was measured by the 18-item Mini-Nutritional Assessment (MNA). The MNA contained a substantial component of anthropometric measurements as well as subscales for dietary behavior, general assessment and subjective health. Each item in MNA can have two or three or four possible choices (Guigoz *et al.*1996).

The Comorbidity Checklist was used to assess the presence of 14 chronic illnesses, including myocardial infarction, angina, heart failures, other heart disease, hypertension, diabetes, arthritis, stroke, lung disease, vision prob-

lems, hearing problems, Parkinson's disease, hip fracture, and cancer(Guralnik ,1989). Each item has two choices(yes=1,no=0), and we sum up the scores of the 15 items of Comorbidity Checklist. SumC is denoted the total score of Comorbidity Checklist. Polypharmacy, operationally defined as the number of medications taken, was assessed by a medication review. Subjects were asked whether, currently, they have taken any prescriptive or over-the-counter medications and, if so, to show the interview all these medications. The total number of medications was treated as a continuous variable.

Depressive symptoms were measured by the 30-item Geriatric Depression Scale (GDS). And each item in GDS are two possible choices (yes=1, no=0) for participants. Here, SumD is used to be the total score for GDS. Using a cutoff score of 11or above, the scale is 84% sensitive and 95% specific for diagnosing depression in the elderly (Yesavage *et al.*,1983).

Functional status, was measured by the 10-item Enforced Social Dependency Scale (ESDS). The ESDS measures physical and social competence. Physical competency includes six activities: eating, dressing, walking, traveling, bathing, and toileting. Social competence includes home, work, and recreational activities, and communication. SumE is denoted the sum score of ESDS, the sum scores range from 10 to 51, with higher scores reflecting greater dependency.

Satisfaction with support was measured by the subscale of Social Support Questionnaire Short Form (SSQSF). In 6 common situations, subjects were asked to list up to nine people who could be counted on (number score) and specified overall degree of satisfaction (satisfaction score). There are

12 items in SSQSF, and each items can be six or nine possible choices for elder to ask. Here, SumS is the total score of SSQSF items.

### 3.2.3    Statistical methods

Data were analyzed using Mplus, version3.13. We used factor analysis to construct the underlying structures of MNA and GOHAI. Factor analysis seeks to discover the observed items of MNA and GOHAI to be explained largely in terms of a much smaller number of factors. Typical factor analysis is for continuous measured variables. Since MNA and GOHAI items are categorical, categorical version of factor analysis is used. For categorical observed variables, the liability threshold model is used, which postulates the existence of an unobserved continuous variable and a set of thresholds. There are two types of factor analysis, we first used exploratory factor analysis (EFA) to determine the number of factors and construct hypothetical categorization of measured items. Generally speaking, we should like to see a pattern of loadings such that each item loads highly on a single factor and has small-to-moderate loadings on the remaining factors. Rotated loadings can evaluate to find a meaningful interpretation of the original data. Here, we did promax rotation which is assumed to be correlated among the factors. We then used confirmatory factor analysis (CFA) to verify the categorization determined by EFA.

The factor analysis for MNA and GOHAI has three stages. At the first stage, we use EFA to determine the number of factors that are needed to explain the correlations among observed MNA and GOHAI items, respec-

tively. After deciding the number of factors, the second stage used CFA to specify the relationships among factors and determine how the factors will be measured. At final stage, structural equation modeling (SEM) was used to build up the relationships between MNA factors and GOHAI factors. We also add some covariates (Age, Gender, Education, SumC, Number of medications, SumD, SumE and SumS) into SEM model to adjust for possible confounding effects for the relation between nutrition and oral health.

In Mplus, it provided some fit indices: weighted root mean square residual (WRMR), the Tucker-Lewis Index(TLI), the Comparative Fit Index(CFI), and the root mean square error of approximation (RMSEA). Hu and Bentler (1999) suggests the following fit index cutoff value guide for good models with continuous outcomes: $TLI > 0.95$, $CFI > 0.95$, $RMSEA < 0.06$. Simulation studies in Yu and Muthén (2001) suggest that these cut off values are reasonable also for categorical outcomes. Yu and Muthén (2001) suggests $WRMR < 0.90$ for good models with continuous as well as with categorical outcomes. We used these criteria to judge whether a good model fit was obtained.

## 3.3 Results

### 3.3.1 The descriptive analysis

Two hundred forty three elders participated in the questionnaire. They were able to provide complete data on all variables of interest and to answer

questions. Participants' mean standard deviation age was $81.6 \pm 9.4$. The characteristics of the study population are reported in Table A.1. The mean scores standard deviation of SumC, SumE and SumS were $4.2 \pm 1.8$, $20.2 \pm 6.4$ and $46 \pm 10.7$. For SumD used a cutoff score of 11 or above, 89.7% reported scores below 11. Initially, there were eighteen observed indictors in MNA. Due to one or more zero cells, so we deleted three variables: independent living (N7), pressure sore (N9), mid-Arm circumference (N17) from the nutrition. If the proportion of a category of a variable was below 0.025, this category was merged with other category. We remerged the some variables as the follows: intake decline (N1) merged "sever loss of appetite"(0) and "moderate loss of appetite"(1), depression/dementia (N5) merged "sever dementia or depression"(0) and "mild dementia"(1), number of meals/day (N10) merged " one meal"(0) and "two meals"(1), self view of nutrition (N15) merged "view self as being malnourished"(0) and "is uncertain of nutritional state"(1); fluid intake (N13) merged "intake at least one serving of dairy products"(0) and "intake two or more servings of legumes"( 0.5).

In general, the correlation matrix was used as the input for the measurement and structural model testing. Hence, spearman correlation is calculated to examine the correlation between each pair of oral health and nutritional outcomes, respectively. The correlation coefficients of spearman correlation test are 0.77 for O1 and O2, 0.75 for O1 and O5, 0.74 for O2 and O5, 0.54 for O2 and O9, 0.5 for O6 and O11, 0.54 for O9 and O10, 0.5 for N6 and N18, the results have shown strong positive correlations between

two observed indicators (see Table A.2 and Table A.3).

## 3.3.2 Exploratory factor analysis and confirmatory factor analysis

At first, we use EFA to determine the number of factors for the items of GOHAI and MNA, respectively. However, determining the optimal number of factors to extract is not a straightforward task since the decision is ultimately subjective. The Kaiser criterion which was proposed by Kaiser in 1960 retains only factors with eigenvalues greater than one and it is probably the one most widely used. In our oral health data, using this criterion, we would retain three factors. In EFA , two-factor fits the data with $\chi^2 = 29.303$ with $df = 19$, $p-value = 0.0614$, and $RMSEA = 0.047$; three-factor fits the data with $\chi^2 = 17.803$ with $df = 16$, $p-value = 0.3355$, and $RMSEA = 0.022$. The chi-square test of model fit is non-significant, indicating that the null hypothesis that the model fits the data cannot be rejected (the model fits the data well). Although two-and three-factor solutions indicated a good fit between data and model ($RMSEA < 0.06$), finally we specified the number of factors to be two (see Figure A.1). Although the three-factor model fits the data is better than two-factor model, there are two observed indicators O6 and O10 which have negative residual variances. Theoretically, such situation suggests that too many factors are being extracted, so we could accept the solution for one less factor, or not use these variables (by Mplus discussions). Therefore, we use two factors

instead of three factors for GOHAI.

Secondly, we perform CFA to confirm and get the significance of factor loadings in EFA, and we would give the names for each factor based on significant loadings. The items significantly loaded on the first factor, Factor1, are limit the food (O1), trouble biting/chewing (O2), swallow comfortably (O3), eat without discomfort (O5), happy with looks of teeth/gums/dentures variable (O7). These items are orally presented and require physical function. Therefore, this factor may be named" Physical functioning of oral health". The second factor, Factor2, is identified by the following items: prevent from speaking (O4), limit contacts with others (O6), use med to relieve pain (O8), worried or concerned with oral problems (O9), feel nervous/self-conscious (O10), uncomfortable eating in public (O11), gum/teeth sensitive to hot/cold/sweets (O12). Most of these items have a common feature in them: these variables are concerned on the self-image, so this factor may be named" Social functioning of oral health"(see Figure A.1 for oral health data). This CFA model fits the oral health data well ($CFI = 0.992$, $TLI = 0.993$, $RMSEA = 0.055$, $WRMR = 0.936$). CFA results also show that the correlation between Factor1 and Factor2 was statistically significant (Figure1). Hence, it could explain that the physical function could affect the social function of oral health.

Moreover, for nutrition data, four-factor model fits to data is better than three-factor model (three-factor fits the data with $\chi^2 = 55.217$ with $df = 36$, $p-value = 0.0212$, and $RMSEA = 0.047$; four-factor fits the data with $\chi^2 = 37.214$ with $df = 29$, $p-value = 0.1408$, and $RMSEA = 0.034$),

but we only take three factors. When fitting four-factor model for MNA, N10 will form a new factor by itself, which cause the factor covariance matrix is not positive definite, thus the variance estimates of factors can not be calculated. Hence, we use three factors instead of four factors for MNA.

In CFA, there are six items including intake decline (N1), weight loss (N2), stress/acute disease (N4), number of meals/day (N10), protein intake (N11), self view of nutrition (N15) in F1(Figure A.1). The first factor, F1, is interpreted as" the nutritional health appraisal." But this item N11 ($p-value = 0.37$) is not significantly loaded on F1. The second factor, F2, has five items including BMI (N6), >3 meds (N8), feeding mode/help with eating (N14), self view of health (N16), calf circumference (N18), it focus on "the general health appraisal." Here, the N8 ($p-value = 0.23$) and N16 ($p-value = 0.85$) items are not significantly loaded on the F2. The final factor, F3, has four items including mobility (N3), depression/dementia (N5), fruit/veggie intake (N12), fluid intake (N13) and it is interpreted as "the dietary behavior appraisal" (see Figure A.1 for nutrition data). This confirmatory three-factor model of MNA excluding N7, N9 and N17 displayed the goodness–of – fit measures as follows: $CFI = 0.876$, $TLI = 0.868$, $RMSEA = 0.056$, $WRMR = 1.079$. As mentioned above, the correlations among those factors, it could show some potential relationships. The measurement model shows the estimate of correlations among factors significantly, as illustrated in Figure A.1. In our sense, the impact of the dietary behavior on the nutritional health is major factor, and the general

30

health has influenced by the nutritional health.

### 3.3.3 The structural equation model

In addition to EFA and CFA, we use Mplus to fit structural equation models that feature causal relationships among latent variables (factors). The Figure A.2 below shows a structural equation model for two independents (Factor1 and Factor2) as causes of three dependents (F1, F2 and F3). It has shown the SEM model fit, most of observed indicators corresponding to their factors are significant except N11, N8 and N16. The physical function (Factor1) influences the general health (F2) ($p - value = 0.005$), and the social function (Factor2) affects two additional factors: the nutritional health (F1) ($p - value = 0.035$) and the general health (F2) ($p - value = 0.015$). Then, Figure A.3 presents the SEM model after adjusting for Age, Gender, Education, SumC, Number of medications, SumD, SumE and SumS. Most of observed indicators in their factors are significant except N11 , N8 , N16 and N5. Unfortunately, there are no significantly casual relationships among those nutrition and oral health factors. However, there were some significant associations between nutritional factors and risk effects as shown in Table A.4, for instance, age affects the nutritional health (F1) and the general health (F2), the influence of gender on the nutritional health (F1) and the general health (F2) are significant. SumD only affect the nutritional health (F1), and SumE will affect the general health (F2) and the dietary behavior factor (F3). The nutritional health (F1) and the dietary behavior (F3) are influenced by SumS.

Comparing two SEM models, Table A.5 presents the fit indices to the data of two models with or without adjusting for risk factors. The chi-square tests of overall model fits are statistically significant for both models, they suggest that the two models may need some modification before they fit the data well. However, both values of RMSEA are well below the recommended 0.06 cutoff that indicates good model fit, so we still consider that these two models fit the data well. It is interesting to find that the causal relationships between nutritional factors and oral health factors disappeared after adjusting for identified risk factors. The confounding effects may be occurred due to SumD, SumE and SumS. Also, Age is important component because older people have nutrition and oral health related problems.

## 3.4   Discussion

1. Discuss the fit indices.

Since the chi-square test is sensitive to sample size (such that large samples often return statistically significant chi-square values) and non-normality in the input variables, Mplus also provides the RMSEA statistic. The RMSEA is not as sensitive to large sample sizes. Sample size has also been shown to be a prominent factor that affects the performance of model fit indices. Because of our sample size, we don't discuss the fit index SRMR, which cut off does not work well with small sample sizes ($N \leq 250$) of this article. A cutoff value close to 1.0 for WRMR is suitable under most

conditions but it is not recommended for latent growth curve models with more time points. CFI performs relatively better than TLI and RMSEA, and a cutoff value close to 0.96 for CFI has acceptable rejection rates across models when $N \geq 250$ (Yu,2002).

2. How to select the model?

For measurement part of oral health, theoretically, O7 (happy with looks of teeth/gums/dentures) should be load on Factor2 (Social functioning of oral health), however, O8 and O12 (gum/teeth sensitive to hot/cold/sweets) should be load on Factor1 (Physical functioning of oral health). After transferring these indicators, the result is not as good as expected ($\chi^2 = 42.905$ with $df = 20$, $p - value = 0.0021$, $CFI = 0.987$, $TLI = 0.989$, $RMSEA = 0.069$, $WRMR = 1.048$). To select suitable model is most difficult step because the computer can't help at this stage, you actually have to think on your own. You specify the model, based on your knowledge of the field, on your reading of the literature, or on theory. Sometimes the sample size results in model unstable.

# Chapter 4

# LONGITUDINAL ANALYSIS FOR MULTIPLE NUTRITIONAL DATA

## 4.1 Introduction

Study nutrition is important because the nutrition is key point for hospitalized elderly patients' general health. Many studies indicated that the cognition, nutrition, and function deteriorate steadily during hospitalization and such significantly affect clinical outcomes for elder people. The purpose of this article is to observe the changes of nutritional status for hospitalized elder patients. Thus, we focus on the nutritional status and some related risk factors.

Here, nutrition is measured by multiple indicators. Previous study uses

the total sum of these indicators to repeat the nutrition status. This might mass the true nutrition. We use factor analysis to draw the distinct dimension of nutrition. Previous studies have collected cross-sectional data of nutritional status, not longitudinal data. However, the main advantage of a longitudinal study is its effectiveness for studying change. Another merit of the longitudinal study is its ability to distinguish the degree of variation in response variables across time for one person from the variation in response variables among people (Digglen *et al.*,2002). Therefore, longitudinal data enable us to study the association between nutrition and risk factors at different stage of hospitalization, and to study the nutrition change and its associated causes.

Due to longitudinal data, in this work we applied generalized estimating equations (GEEs) for the analysis of nutritional changes for hospitalized older people. The advantage of the GEE approach is that it requires weaker distributional assumptions and maintains the properties of consistency and asymptotic normality of parameter estimates. And useful of the GEE approach is that it is not necessary for the "working" correlation matrix to be correctly specified to construct consistency and asymptotic normality of parameter estimates (Zeger and Liang,1986).

In this paper, we combine factor analysis and longitudinal data analysis to study the nutritional changes at different dimension of nutrition and investigate interactions among some risk factors and nutritional factors in section 4.2 and the results for those statistical methods would present in section 4.3. And the section 4.4 would discuss some findings or statistical

methods.

## 4.2 Metohd

### 4.2.1 Sample

Hospitalized patients age 65 years and older of National Taiwan University Hospital would be collected into the study. The sampling frame is 20 medical and surgical adult units (excluding oncology, AIDS, and Hospice units) from a 2500-bed tertiary medical center. Four units (2 medical and 2 surgical units) would be randomly selected as the study sties. Face-to-face assessments would be conducted with standardized measures by trained research nurses in four data collection points. Subjects with severe cognitive impairment would be excluded, since the study design involved use of self-report questions. Scoring less than 20 in the Chinese Mini-Mental State Exam (MMSE) at baseline would meet this exclusion criterion. After completing the baseline assessment before older patients discharge to home, each subject would be followed for six months. Up to February, 2006 , the four data collection points contain patients as follows: 302 patients within 48 hours of admission (Time1), 288 patients before discharge (Time2), 217 patients for three month post index hospitalization (Time3), 167 patients for six month post index hospitalization (Time4).

## 4.2.2 Instruments

Chen *et al.*(submitted) proposed possible factors which impacted on malnutrition in the elderly. In that paper, they provided four conceptual texts including loss, chronic illness, dependency and loneliness to test the relationships with the development of malnutrition in the elderly. In loss part, it contains some empirical indices including age, oral health. The chronic part includes as follows: comorbidity, number of medications. For illness, there are several indices including gender/education, functional status. And the loneliness main contains social support and depression empirical indices.

Demographics and Treatment-related Chart Data

A demographic form will be designed to collect the data including age at admission, gender, marital status, living status, education, and ethnic group, etc. Additionally, chart data tracking form will be developed to elicit the admission diagnosis, LOS, NPO days, cost of care (from the hospital billing system) and laboratory data including serum albumin, cholesterol and count of blood cell data.

Comorbidities

The history and number of comorbidities will be elicited from the medical record. A standardized comorbidity checklist will be used to assess common chronic illnesses including myocardial infarction, angina, heart failure, hypertension, diabetes, Hyperlipidemia, arthritis, stroke, lung disease, renal disease, version problems, hearing problems, Parkinson's dis-

ease, osteoporosis, hip fracture, pressure sore, and cancer, etc. There are 20 items in Comorbidity list, and each item can have two possible choices (yes=1, no=0). SumC is used to be the total score of comorbidity items.

### Medication

Medication review will be conducted. The number and type of prescription and over-the-counter medications taken by subjects will be documented. A protocol will be developed.

### Oral Health

A 12-item Chinese General Oral Health Assessment Index (GOHAI) will be used to assess oral health. The GOHAI is designed to assess the dimensions of oral function (eating and speaking), pain, discomfort, worry, and oral health related social functioning (Atchison & Dolan,1990). Information on number of remaining teeth, dental status(full denture, bridge, partial denture, or natural teeth), fitness of denture, and dental care utilization (regular dental check-up or not) will also be solicited. Each item can have five possible choices (always/often/sometimes/seldom/never), and SumO is the total score of these items in GOHAI.

### Nutritional Status

The 18-item Chinese version of Mini-Nutritional Assessment (MNA) along with some other nutritional makers will be used to measure nutritional status. Additionally, weight loss per unit of time, percentage of usual weight, and serum albumin and cholesterol levels will also be reported. Each item in MNA can have two or three or four possible choices.

### Funtional Status

Functional status, was measured by the 10-item Enforced Social Dependency Scale (ESDS). The ESDS measures physical and social competence. Physical competency includes six activities: eating, dressing, walking, traveling, bathing, and toileting. Social competence includes home, work, and recreational activities, and communication. Each item can have three or four or six possible choices, and SumE is used to count the total score of ESDS.

### Depressive Symptoms

The 30-item Chinese version of Geriatric Depressive Scale (GDS) will be used to measure the presence of depressive symptoms. The Institute of Medicine has recommended GDS for clinical use. Here, we only use 15-items GDS. Each item states depressive symptoms and participants are asked to indicate yes/no they feel the way described. SumD is counted the total score of the 15 items in GDS.

### Social Support

Satisfaction with support was measured by the subscale of Social Support Questionnaire Short Form (SSQSF). In 6 common situations, subjects were asked to list up to nine people who could be counted on (number score) and specified overall degree of satisfaction(satisfaction score). There are 12 items in SSQSF, and each items can be seven or ten possible choices for elder patients to ask. Here, SumS is the total score of SSQSF items.

### 4.2.3 Statistical analysis

Data were analyzed using Mplus,version 3.13 and R package. We measured the nutrition situation at four different time points: within 48 hours of admission (Time1), before discharge (Time2), three month post index hospitalization (Time3), six month post index hospitalization (Time4).

For Time1 nutrition data, do EFA to determine the number of factors and possible indicators of each factor. The relations obtained from Time1 are applied to other time points. In this study, we want to know how the nutrition changed over time and to predict future nutritional status from the preceding ones, thus , it was reasonable to apply the factor structure obtained from Time1 to all other time points. We then combined all the data from four time points. For each factor of all four time points, do CFA as shown in Figure A.4-Figure A.6, and the relations between each indicator and factor at all four time points are all constrained to be equal. We obtained the factor scores for each factor and time points based on the above CFA models.

For each factor, the factor scores from four time points are correlated. We used GEE approach to model the correlations. In the following, we take the Factor1 as an example to illustrate the implemented GEE model. Let Factor11, Factor12, Factor13 and Factor14 be the factor scores of Factor1 from four time points, respectively. The primary variables of interest in these analyses were the indicators of Time2($t_2$), Time3($t_3$ ) and Time4($t_4$), Age, Gender, Education($EDU$), SumO, SumC, Number-

of-medication($NM$), SumD , SumE and SumS. The interaction terms of
the time point indicators with other variables were also put into the model.
The marginal model for factor scores of Factor1, $Factor1score$ , as a function of the covariates is assuming to be followings:

$$
\begin{aligned}
E(Factor1score) \;=\; & \beta_0 + \beta_1 t_2 + \beta_2 t_3 + \beta_3 t_4 + \beta_4 Age + \beta_5 Gender \\
& + \beta_6 EDU + \beta_7 SumO + \beta_8 SumC + \beta_9 NM + \beta_{10} SumD \\
& + \beta_{11} SumE + \beta_{12} SumS + \beta_{13} t_2 Age + \beta_{14} t_3 Age \\
& + \beta_{15} t_4 Age + \beta_{16} t_2 Gender + \beta_{17} t_3 Gender \\
& + \beta_{18} t_4 Gender + \beta_{19} t_2 EDU + \beta_{20} t_3 EDU + \beta_{21} t_4 EDU \\
& + \beta_{22} t_2 SumO + \beta_{23} t_3 SumO + \beta_{24} t_4 SumO + \beta_{25} t_2 SumC \\
& + \beta_{26} t_3 SumC + \beta_{27} t_4 SumC + \beta_{28} t_2 NM \\
& + \beta_{29} t_3 NM + \beta_{30} t_4 NM + \beta_{31} t_2 SumD + \beta_{32} t_3 SumD \\
& + \beta_{33} t_4 SumD + \beta_{34} t_2 SumE + \beta_{35} t_3 SumE + \beta_{36} t_4 SumE \\
& + \beta_{37} t_2 SumS + \beta_{38} t_3 SumS + \beta_{39} t_4 SumS \quad\quad (4.1)
\end{aligned}
$$

where $t_2= 1$ if measures at Time2 and 0 if not,

$t_3= 1$ if measures at Time3 and 0 if not,

$t_4= 1$ if measures at Time4 and 0 if not.

From this model, we can obtain the relationship between factor scores
and risk factors at each time point, and the relation of the nutritional
change between two adjacent time points with risk factors. The relationship
between Factor1 and, for example, Age at Time1 can be represented by the

regression coefficient $\beta_4$ , which is the Factor1 score change per one-year increase in age ; at Time2 is $\beta_4 + \beta_{13}$; at Time3 is $\beta_4 + \beta_{14}$ ; and at Time4 is $\beta_4 + \beta_{15}$. Also, let's take Age as an example to show how the difference in factor scores between two adjacent time points related to risk factors. For every one year increase in age, the difference in Factor1 score between Time2 and Time1 has a $\beta_{13}$-unit change. Therefore, if $\beta_{13}$ is positive, it means that the older the people, the larger the difference in Factor1 score between Time2 and Time1. Furthermore, the difference in Factor1 score between Time3 and Time2 has a $(\beta_{14} - \beta_{13})$-unit change and the difference between Time4 and Time3 has a $(\beta_{15} - \beta_{14})$- unit change, for every one-year increase in age.

The model for the correlation among Factor1 scores in four time points is modeled as

$$corr(Factor1score_j, Factor1score_t) = \alpha \quad for \quad j < t = 1, 2, 3, 4.$$

We here assumed an "exchangeable" correlation model.

## 4.3   Results

### 4.3.1   The descriptive analysis

Hospitalized elders' age ranged from 64 to 89 years with mean of 72.14 years and standard deviation of 5.72 years. The characteristics of the study population is reported in Table A.6. At Time1, the mean scores $\pm$ standard

deviation of SumO, SumC, Number of medications, SumD ,SumE, and SumS were $50.30 \pm 6.57$, $15.86 \pm 1.80$, $3.27 \pm 2.57$, $4.58 \pm 3.52$, $16.11 \pm 6.26$ and $53.52 \pm 10.44$, respectively ; at Time2 were $49.62 \pm 6.36$, $15.73 \pm 1.66$, $4.22 \pm 2.01$, $7.56 \pm 3.73$, $29.79 \pm 7.91$ and $51.29 \pm 10.70$, respectively ; at Time3 were $47.86 \pm 5.93$, $15.12 \pm 2.59$, $4.01 \pm 1.62$, $5.21 \pm 3.97$, $20.47 \pm 8.13$ and $52.14 \pm 10.06$, respectively ; at Time4 were $47.33 \pm 6.02$, $15.51 \pm 2.30$, $3.73 \pm 1.67$, $5.32 \pm 4.15$, $21.18 \pm 7.65$ and $50.48 \pm 9.27$, respectively. Above the mean scores , SumO changed decreasingly over time. The mean scores of others also changed over time excluding SumC. They neither increased nor decreased strictly over time.

Initially, there were eighteen variable indictors in MNA. Due to one or more zero, so we deleted three variables psychological stress/recent illness (N7), depression/dementia (N9), pressure sore (N10) from the nutrition data. If the proportion of a category of a variable was below 0.025, this category was merged with other category. We remerged the some variables as the follows: weight loss (N4) in Time3 merged "weight loss greater than 3kg"(1) and "weight loss between1 and 3kg"(2); mobility (N8) for all Time1 to Time4 merged "bed or chair bound"(0) and "able to get out of bed/chair but does not go out"(1); number of meals/day (N11) for all Time1 to Time4 merged "one meal"(0) and "two meals"(1); feeding mode/ help with eating (N16) for all Time1 to Time4 merged "unable to eat without assistance"(0) and "self-fed with some difficulty "(1).

### 4.3.2 Exploratory factor analysis at Time1

We used Time1 to construct the relations between indicators and factors, and then the results are applied to other time points. At first, we need to decide the number of factors by doing EFA in Mplus, the EFA fit measure as the follows: $\chi^2 = 55.724$ with $df = 39$, $p-value = 0.0402$, $RMSEA = 0.038$, hence we select three factors at Time1. We took two different categorization of measured indicators to perform CFA and compare the fit indices between the two (see Table A.7). The first categorization which is in line with EFA, the results has shown as follows: there are three items including N1, N2 and N3 in Factor1. The second factor, Factor2, has five items including N4, N11, N12, N14 and N17 .The final factor, Factor3, has seven items including N5, N6, N8, N13, N15, N16 and N18. However, in the second categorization, we put the observed indicators N13 and N15 into Factor2 because these two indicators are thought to be more similar in theory to indicators of Factor2. The CFA fit has shown in Table A.7. Both the values of chi-square are significant, and it means that two CFA models did not fit the data well. The RMSEA shows that the second categorization is better than the first categorization. Finally, we still took the first categorization for analyzing the data. Since the residual variance for N11 variable was not positive from the second categorization, so the factor scores of Factor2 at Time3 could not be computed. Hence, we could not use the second categorization.

### 4.3.3 Confirmatory factor analysis combining four time points

The categorization of indicators at Time1 was applied to other time points. We then do CFA for each factor combining all four time points. The relation between each indicator and the factor are constrained to be equal across different time points to make sure the measurement invariance of the factors across time. After obtaining the model, the Figure A.4-Figure A.6 displayed the relations among indicators and each factor. Factor1 (Figure A.4) is interpreted as " Anthropometric indices." When the score of mid arm circumference (N2) is high, it would decrease in Factor1. But the higher score of calf circumference (N3), it would increase in Factor1. That is, the higher scores of N2 and N3 indicators mean that the person is fatter than others. Factor2 (Figure A.5) is interpreted as "Nutritional behavior indices." When the patients get high scores for those items as follows: number of meals/day (N11), protein intake (N12), intake decline (N14), self view of nutrition (N17), it would increase in Factor2. That is, their nutritional behaviors are better than others. Factor3 (Figure A.6) is interpreted as "Nutritional risk indices. If the patients get high scores as follows: >meds(N6), mobility(N8), fruit/veggie intake(N13), fluid intake(N15), feeding mode/help with eating(N16), it would increase in Factor3.

## 4.3.4  Generalized estimating equations(GEEs)

The results of the GEE analysis predicting risk factors which influence nutrition at each time point are reported in Table A.8. For Factor1 (Anthropometric indices), only Age and SumD are significant at Time1. For Factor2 (Nutritional behavior indices), SumO and SumD are significant over different time. For Factor3 (Nutritional risk indices), Number of medications, SumD and SumE decrease significantly over different time. Hence, we could make other comments. We summarized the results from the Table A.8. At Time1 (with 48 hours of admission), each nutritional factor (Facto1 - Factor3) is influenced by SumD, thus depressive symptoms would affect the nutritional status. Besides, Factor3 is influenced by most risk factors. For Time2 (before discharge before), some risk factors are statistical significant especially Gender, SumD and SumE. For Time3 (three month post index hospitalization), the impacts of SumD and SumE on Factor2 and Factor3 are major factors. For Time4 (six month post index hospitalization), SumO and SumD would affect nutritional status significantly.

Table A.9 demonstrated the results GEE for associations between nutritional difference and identified risk factors. The nutritional difference is defined as Time1 to Time2, Time2 to Time3, Time3 to Time4. An interaction term of time by risk factors that is statistically significant with p-value <0.05 indicates that the risk factor is associated with the difference of nutritional status. For instance, the difference in Factor1 (Anthropometric indices) score between Tim1 and Time2 are significant associated by Num-

ber of medications and SumE. And the difference in Factor2 (Nutritional behavior indices) score between Time1 and Time2 has largely changed by SumC and SumD. Some risk factors would affect the differences in Factor 3 score between Time1 and Time2 or Time2 and Time3 such as Gender, Education, Number of medications, SumE. Moreover, the difference in Factor3 score between Tim3 and Time4 has largely changed by SumO, Number of medications and SumS.

## 4.4　Discussion

For longitudinal data, there are two distinct types: subject-specific (SS) models and population-averaged (PA) (Zeger,Liang&Albert,1988).To model SS, the random effects model can be used (Zeger and Liang,1992). In general, it is useful for doctors to do individual patients diagnosis. However, the GEE approach used in this thesis is for modeling PA. We here focused on PA models which describe how the average response across subjects changes with covariates. For epidemiology, we applied the average of the population to understand all characteristics for the population.

# Chapter 5

# CONCLUSION

Of this thesis, it contains two primary sections. First, we are interested in discovering the relation between nutrition and oral health. Since nutrition and oral health are multiple indicators, the former studies used total sums of MNA and GOHAI to do analysis. To get the true association, we use factor analysis to draw distinct dimension of nutrition. From the results, we took two SEM models with or without adjusting for identified risk factors. Those risk factors are used to adjust for possible confounding effects for the relation between nutrition and oral health. Besides, Mplus provided the some useful indices for readers to examine how the model fits to data. Here, we mainly use the RMSEA value with cutoff 0.06 to evaluate CFA and SEM models. However, some of chi-square values for models are statistically significant, then we would reject null hypothesis. It means the models don't fit to data well. Maybe we have to collect large sample size in the further study.

In the second part, we take longitudinal data for nutrition instead of cross-sectional data. The purpose of this section is to observe the relationships of nutritional factors with risk variables at different time points and understand how the difference in factor scores between two time points related to risk factors. Thus, we here apply factor analysis and generalized estimating equations (GEEs) approach. Also, from the results, SumD played an important role in the relationship between factor score at each time point. However, it was not significant for the difference in factor score between two adjacent time points. On the contrary, SumO had reverse conditions. Education seemed to have no association with factor score at all four time points. But it was significant for the difference in Factor3 (Nutritional risk indices) scores between Time1 and Time2. Similarly, SumS was only significant for the difference in Factor3 scores between Time3 and Time4. To sum up, most of the identified risk factors would affect nutritional status. It may improve the nutritional status with controlling those identified risk factors for hospitalized elderly patients.

# References

[1] Atchison,K.A.and Dolan,T.A.(1990):"Development of the geriatric oral health assessment index," *Journal of Dental Education*,**54**,680-687.

[2] Basilevsky,A. (1994): Statistical factor analysis and related methods.

[3] Bartlett,M.S.(1937),"The Statistical Conception of Mental Factors," *British Journal of Psycholog*,**28**,97-104

[4] Bollen,K.A.(1989):Structural equations with latent variables.

[5] Browne,M.W.and Cudeck,R.(1993): Alternative ways of assessing model fit. In K.Bollen and J.S.Long(eds.),Testing Structural Equation Models (pp.136-162). Newbury Park:Sage Publications.

[6] Chen,C.-H.;Bai,Y.-Y.;Huang,G.-H.and Tang,S.T.:"Revisiting the Concept of Malnutrition in the Elderly."(submitted).

[7] Diggle,P.J.;Heagerty,P.J.;Liang,K.-Y.and Zeger,S.L.(2002): Analysis of Longitudinal Data.Oxford University Press, New York, second edition.

[8] Fitzmaurice,G.M. and Laird,N.M.(1993):"A likelihood-based method for analyzing longitudinal binary responses ,"*Biometrika*,**80**,141-151.

[9] Fitzmaurice,G.M.;Laird,N.M.and Rotnitzky,A.G.(1993):"Regression Models for Discrete Longitudinal Responses,"*Statistical Science*,**8**,284-299.

[10] Guralnik,J.M.(1989): Aging in the eighties:the prevalence of comorbidy and its a ssociation with disability.Advance data from vital and health statistics,170.Hyattsiville,MD:National Center for Health Statistics.

[11] Guigoz,Y.;Vellas,B.and Garry,P.J.(1996):"Assessing the nutritional status of the elderly:the mini nutritional assessment as part of the geriatric evaluation," *Nutrition Review,* **54**,S59-S65.

[12] Hu,L.T.and Bentler,P.M.(1999):"Cutoff criteria for fit indices in covariance structure analysis:conventional criteria versus new alternatives," *Structural Equation Modeling*,**6**,1-55.

[13] Horton,N.J.and Lipsitz,S.R.(1999):"Review of Software to Fit Generalized Estimating Equation Regression Models." *The American Statistician*,**53**,160-169.

[14] Liang, K.-Y.and Zeger, S.L.(1986):"Longitudinal Data Analysis Using Generalized Linear Models,"*Biometrika*,**73**,13-22.

[15] Liang,K.-Y.,Zeger,S.L.and Qaqish,B.(1992):"Multivariate regression analyses for categorical data(with Discussion)," *Journal of the Royal Statistical Society*, **B**,**54**,3-40.

[16] Lee ,Y.and Nelder,J.A.(1996):"Hierarchical Generalized Linear Models," *Journal of the Royal Statistical Society B*,**58**(4),619-678.

[17] Lee ,Y.and Nelder,J.A.(2001):"Hierarchical Generalized Linear Models:A Synthesis of Generalized Linear Models,Random-Effect Models and Structured Dispersions." *Biometrika*,**88**(4),987-1006.

[18] Muthén, B.(1979):"A structural probit model with latent variables."*Journal of the American Statistical Association*,**74**,807-811.

[19] Muthén, B.(1981a),A general structural equation model with ordered categorical and continuous latent variable indicators(Department of Psychology, University of California, Loss Angeles, CA)

[20] McCullagh, P.(1983):"Quasi-likekihood functions," *Annals of Statistics*,**11**,59-67

[21] McCullagh, P.and Nelder,J.A.(1983):Generalized Linear Models. London:Chapman and Hall.

[22] Muthén,B.(1983):" Latent variable structural equation modeling with categorical data," *Journal of Econometrics*, **22**, 48-65.

[23] Muthén,B.(1984):"A general structural equation model with dichoto-
mous,ordered categorical, and continuous latent variable indica-
tors,"*Psychometrika*, **49**, 115-132.

[24] Muthén,B.(1989b):"Latent variable modeling in heterogeneous popu-
lations,"*Psychometrika*,**54**,557-585.

[25] Muthén, B.,and Kaplan,D.(1992):"A comparison of some methodolo-
gies for the factor analysis of non-normal likert variables," *British Jour-
nal of Mathematical and Statistical Psychology*,**38**,171-189.

[26] Mueller, R.O.(1996): Basic principles of structural equation modeling.

[27] Nicholas,J.H. and Stuart, R.L.(1999):" Review of Software to Fit Gen-
eralized Estimating Equation Regression Models, " *The American Sta-
tistician*,**53**,160-169.

[28] Prentice,R.L.(1988):"Correlated binary regression with covariates spe-
cific to each binary observation," *Biometrics*,**44**,1033-1048.

[29] Reyment,R.and Joreskog,K.G.(1993): Applied factor analysis in the
natural sciences.

[30] Ritchie, C.S.; Joshipura K.; Hung, H.C.; Douglass,C.W.(2002): "Nu-
trition as a mediator in the relation between oral and systemic disease:
Associations between specific measures of adult oral health and nutri-
tion outcomes."

[31] Skrondal,A.and Sophia R.-H.(2005):"Structural equation modeling:categorical variables."

[32] Thall,P.F.and Vail,S.C.(1990):"Some covariance models for longitudinal count data with overdispersion,"*Biometrics*,**46**,657-671.

[33] Wedderburn,R.W.M(1974)."Quasilikelihood function, generalized linear models and Gauss-Newton method,"*Biometrika*,**61**,439-447.

[34] Yu,C.-Y. and Muthén, B.(2001):Evaluation of model fit indices for latent variable models with categorical and continuous outcomes. Technical report in preparation.Newbury Park: SAGE Publications.

[35] Yu,C.-Y.(2002): Evaluating cutoff criteria of model fit indices for latent variable models with binary and continuous outcomes.

[36] Yesavage,J.A.;Brink,T.L.; Rose,T.L.(1983):"Development and validation of a geriatric depression scale,"*Joural of Psychiatric Research*,**17**,31-49.

[37] Zeger,S.L.and Liang, K.-Y. (1986): " Longitudinal data analysis for discrete and continuous outcomes, " *Biometrics*, **42**,121-130.

[38] Zeger,S.L.;Liang,K.-Y.;Albert, P.S.(1988): " Models for Longitudinal Data: A Generalized Estimating Equation," *Biometrics*,**44**,1049-1060.

[39] Zhao,L.P.and Prentice, R.L.(1990):"Correlated binary regression using a generalized quadratic model,"*Biometrika*,**77**,642-648.

[40] Zeger,S.L.and Liang, K.-Y.(1992): " An overview of methods for the analysis of longitudinal data," *Statistics in medicine*, **11**,1825-1839.

[41] Ziegler A.;Kastner C. and Blettner M.(1998):"The Generalized Estimating Equations:An Annotated Bibliography," *Biometrical Journal*,**40**,115-139

# Appendix A

# List of Tables and Figures

Table A.1: Characteristics of community-dwelling elders in U.S.A.

|  | Number | % |
|---|---|---|
| **Demographics** | | |
| **Age Group** | | |
| 74 years and below | 42 | 17.28% |
| 75 to 84 years | 80 | 32.92% |
| 85 to 90 years | 68 | 27.98% |
| 91 years and above | 53 | 21.81% |
| **Gender** | | |
| Male | 52 | 21.4% |
| Female | 191 | 78.6% |
| **Ethnicity** | | |
| White | 168 | 69.14% |
| Other | 75 | 30.86% |
| **Living Status** | | |
| Live Alone | 223 | 91.77% |
| Live with Others | 20 | 8.23% |
| **Marital Status** | | |
| Windows | 159 | 65.43% |
| Married | 19 | 7.82% |
| Divorced | 39 | 16.05% |
| Single/Never Married | 26 | 10.70% |
| **Education** | | |
| Less than high school and hingh school | 160 | 65.84% |
| Greater | 83 | 34.16% |
| **Religion** | | |
| Catholic | 29 | 11.93% |
| Protestant | 73 | 30.04% |
| Jewish | 136 | 55.97% |
| Buddhist | 1 | 0.44% |
| Others | 4 | 1.65% |

Table A.2: Spearman correlation among observed indicators of oral health

| | O1 | O2 | O3 | O4 | O5 | O6 | O7 | O8 | O9 | O10 | O11 | O12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| O1 | 1.0000 | | | | | | | | | | | |
| O2 | 0.7697* | 1.0000 | | | | | | | | | | |
| O3 | 0.2855 | 0.3824 | 1.0000 | | | | | | | | | |
| O4 | 0.2986 | 0.2727 | 0.1113 | 1.0000 | | | | | | | | |
| O5 | 0.7519* | 0.7431* | 0.3722 | 0.3526 | 1.0000 | | | | | | | |
| O6 | 0.2260 | 0.1761 | 0.0244 | 0.4527 | 0.2241 | 1.0000 | | | | | | |
| O7 | 0.4607 | 0.4318 | 0.1655 | 0.2628 | 0.4264 | 0.1363 | 1.0000 | | | | | |
| O8 | 0.1067 | 0.0962 | 0.0693 | 0.2315 | 0.1316 | 0.3335 | 0.0697 | 1.0000 | | | | |
| O9 | 0.4900 | 0.5392* | 0.1906 | 0.2455 | 0.4550 | 0.2943 | 0.3974 | 0.1440 | 1.0000 | | | |
| O10 | 0.3465 | 0.3604 | 0.1940 | 0.2790 | 0.3544 | 0.4794 | 0.3140 | 0.2314 | 0.5412* | 1.0000 | | |
| O11 | 0.3008 | 0.3315 | 0.0245 | 0.2935 | 0.2714 | 0.5009* | 0.1994 | 0.2615 | 0.3321 | 0.4758 | 1.0000 | |
| O12 | 0.1696 | 0.0874 | 0.0822 | 0.1097 | 0.1695 | 0.1277 | -0.0345 | 0.2243 | 0.1002 | 0.0003 | 0.1447 | 1.0000 |

Table A.3: Spearman correlation among observed indicators of nutrition

| | N1 | N2 | N3 | N4 | N5 | N6 | N8 | N10 | N11 | N12 | N13 | N14 | N15 | N16 | N18 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| N1 | 1.0000 | | | | | | | | | | | | | | |
| N2 | 0.2200 | 1.0000 | | | | | | | | | | | | | |
| N3 | 0.1370 | 0.1030 | 1.0000 | | | | | | | | | | | | |
| N4 | 0.2395 | 0.1548 | 0.0737 | 1.0000 | | | | | | | | | | | |
| N5 | 0.1792 | 0.1285 | 0.1448 | 0.0622 | 1.0000 | | | | | | | | | | |
| N6 | 0.2804 | 0.1931 | 0.1100 | 0.1225 | 0.1330 | 1.0000 | | | | | | | | | |
| N8 | -0.1055 | 0.0006 | -0.0057 | 0.1142 | -0.0297 | -0.2342 | 1.0000 | | | | | | | | |
| N10 | 0.1868 | 0.1995 | -0.1511 | 0.0043 | 0.0402 | 0.0189 | -0.0286 | 1.0000 | | | | | | | |
| N11 | 0.0419 | 0.0105 | 0.1351 | -0.0169 | -0.0114 | 0.0093 | 0.0016 | 0.0598 | 1.0000 | | | | | | |
| N12 | 0.1305 | 0.1414 | 0.0355 | 0.0778 | 0.1738 | 0.1444 | 0.0674 | 0.0881 | -0.0376 | 1.0000 | | | | | |
| N13 | 0.0311 | 0.0239 | 0.1213 | -0.0100 | 0.1082 | 0.2675 | -0.0274 | 0.0347 | 0.0118 | 0.2531 | 1.0000 | | | | |
| N14 | -0.1190 | 0.0957 | 0.2161 | 0.0512 | 0.1754 | 0.0710 | 0.0891 | -0.0234 | 0.0783 | 0.0372 | 0.1391 | 1.0000 | | | |
| N15 | 0.4634 | 0.1488 | 0.1528 | 0.3127 | 0.1250 | 0.2120 | 0.0402 | 0.1996 | -0.0227 | 0.1732 | 0.0680 | 0.1787 | 1.0000 | | |
| N16 | -0.0320 | 0.0107 | 0.0078 | 0.1790 | -0.0621 | -0.0467 | 0.1066 | -0.0704 | -0.0650 | 0.1258 | -0.0059 | 0.0350 | 0.0903 | 1.0000 | |
| N18 | 0.2532 | 0.1381 | 0.1308 | 0.1102 | 0.1613 | 0.5062* | -0.0545 | -0.0061 | -0.0466 | 0.0741 | 0.1603 | 0.1114 | 0.1787 | -0.0638 | 1.0000 |

* Results present strong associations between variables

Table A.4: The relationship between each nutritional factor and adjusted risk factors from the SEM model

| Variables | Nutritional health (F1) p-value | General health (F2) p-value | Dietary behavior (F3) p-value |
|---|---|---|---|
| Age | 0.0468* | 0.0080* | 0.1188 |
| Gender | 0.0033* | 0.0280* | 0.4827 |
| Education | 0.6803 | 0.3794 | 0.1352 |
| SumC | 0.1031 | 0.1081 | 0.2325 |
| #Medication | 0.8243 | 0.3865 | 0.0733 |
| SumD | 0.0149* | 0.2117 | 0.1585 |
| SumE | 0.2801 | 0.0407* | 0.0466* |
| SumS | 0.0370* | 0.2960 | 0.0207* |

* Results are statistically significant with p-value<0.05

Table A.5: Goodness-of-fit indices for the SEM models with or without adjusting for risk factors

| | | Unadjusted | Adjusted |
|---|---|---|---|
| **Chi-Square Test of Model Fit** | | | |
| | Value | 108.936 | 145.654 |
| | Degrees of Freedom | 59 | 92 |
| | P-value | 0.0001 | 0.0003 |
| **Comparative Fit Index (CFI)** | | 0.959 | 0.957 |
| **Tucker-Lewis Index (TLI)** | | 0.961 | 0.955 |
| **Root Mean Square Error Of Approximation** | | | |
| **(RMSEA)** | Estimate | 0.059 | 0.049 |
| **Weighted Root Mean Square Residual** | | | |
| **(WRMR)** | Value | 1.137 | 1.120 |

Table A.6: Baseline characteristics of elderly hospitalized patients in National Taiwan University Hospital, March , 2006

|  | Number | % |
|---|---|---|
| **Demographics** | | |
| **Age Group** | | |
| 65 to 74 years | 199 | 65.89% |
| 75 to 84 years | 94 | 31.13% |
| 85 years and above | 9 | 2.98% |
| **Gender** | | |
| Male | 170 | 56.29% |
| Female | 132 | 43.71% |
| **Ethnicity** | | |
| Taiwanese | 242 | 80.13% |
| Others | 60 | 19.87% |
| **Living Status** | | |
| Live Alone | 11 | 3.64% |
| Live with Others | 291 | 96.36% |
| **Marital Status** | | |
| Windows | 81 | 26.82% |
| Married | 214 | 70.86% |
| Divorced | 3 | 0.99% |
| Single/Never Married | 4 | 1.32% |
| **Education** | | |
| Less than high school and hingh school | 255 | 84.44% |
| Greater | 47 | 15.56% |
| **Religion** | | |
| Buddhist | 187 | 61.92% |
| Daoism | 44 | 14.57% |
| Others | 71 | 23.51% |

Table A.7: Two CFA models for nutrition at Time1

| | | 1st categorization | 2nd categorization |
|---|---|---|---|
| **Chi-Square Test of Model Fit** | | | |
| | Value | 95.180 | 86.490 |
| | Degrees of Freedom | 45 | 45 |
| | P-value | 0.0000 | 0.0002 |
| **Comparative Fit Index (CFI)** | | 0.851 | 0.876 |
| **Tucker-Lewis Index (TLI)** | | 0.874 | 0.896 |
| **Root Mean Square Error Of Approximation** | | | |
| **(RMSEA)** | Estimate | 0.061 | 0.055 |
| **Weighted Root Mean Square Residual** | | | |
| **(WRMR)** | Value | 1.197 | 1.133 |

61

Table A.8: The relationship between factor score and risk factors at each time point

| | Time1 R.C. | Time1 95%C.I.+ | Time2 R.C. | Time2 95%C.I.+ | Time3 R.C. | Time3 95%C.I.+ | Time4 R.C. | Time4 95%C.I.+ |
|---|---|---|---|---|---|---|---|---|
| **Factor1 ~ covariates** | | | | | | | | |
| **Age** | -0.0138 | (-0.0261 , -0.0015)* | -0.0108 | (-0.0260 , 0.0044) | -0.0144 | (-0.0315 , 0.0027) | -0.0155 | (-0.0330 , 0.0020) |
| **Gender** | -0.1176 | (-0.2604 , 0.0251) | -0.1289 | (-0.2978 , 0.0399) | -0.0741 | (-0.2845 , 0.1363) | 0.0063 | (-0.2050 , 0.2175) |
| **Education** | -0.0347 | (-0.2184 , 0.1490) | -0.1653 | (-0.3999 , 0.0693) | -0.1280 | (-0.3769 , 0.1209) | 0.0036 | (-0.2764 , 0.2835) |
| **SumO** | 0.0031 | (-0.0041 , 0.0102) | 0.0058 | (-0.0061 , 0.0177) | 0.0078 | (-0.0058 , 0.0213) | 0.0184 | (0.0032 , 0.0336)* |
| **SumC** | 0.0070 | (-0.0172 , 0.0312) | -0.0100 | (-0.0491 , 0.0291) | -0.0089 | (-0.0561 , 0.0383) | -0.0091 | (-0.0568 , 0.0386) |
| **#Medication** | -0.0026 | (-0.0175 , 0.0123) | 0.0175 | (-0.0118 , 0.0467) | -0.0047 | (-0.0509 , 0.0415) | -0.0329 | (-0.0710 , 0.0053) |
| **SumD** | -0.0169 | (-0.0285 , -0.0053)* | -0.0144 | (-0.0366 , 0.0079) | -0.0137 | (-0.0419 , 0.0145) | -0.0173 | (-0.0448 , 0.0103) |
| **SumE** | 0.0053 | (-0.0020 , 0.0126) | -0.0031 | (-0.0177 , 0.0115) | -0.0019 | (-0.0209 , 0.0171) | 0.0001 | (-0.0164 , 0.0165) |
| **SumS** | 0.0024 | (-0.0014 , 0.0062) | 0.0027 | (-0.0047 , 0.0101) | 0.0014 | (-0.0070 , 0.0099) | 0.0032 | (-0.0057 , 0.0120) |
| **Factor2 ~ covariates** | | | | | | | | |
| **Age** | -0.0062 | (-0.0129 , 0.0005) | -0.0056 | (-0.0124 , 0.0012) | 0.0032 | (-0.0027 , 0.0091) | -0.0001 | (-0.0057 , 0.0054) |
| **Gender** | 0.0415 | (-0.0349 , 0.1179) | 0.1070 | (0.0012 , 0.1770)* | 0.0299 | (-0.0345 , 0.0942) | 0.0116 | (-0.0551 , 0.0782) |
| **Education** | 0.0776 | (-0.0071 , 0.1623) | 0.0054 | (-0.0863 , 0.0972) | 0.0324 | (-0.1159 , 0.0511) | 0.0148 | (-0.0698 , 0.0994) |
| **SumO** | 0.0074 | (0.0017 , 0.0130)* | 0.0100 | (0.0042 , 0.0159)* | 0.0106 | (0.0044 , 0.0168)* | 0.0088 | (0.0010 , 0.0167)* |
| **SumC** | 0.0169 | (-0.0038 , 0.0376) | -0.0204 | (-0.0441 , 0.0032) | 0.0004 | (-0.0119 , 0.0128) | -0.0029 | (-0.0160 , 0.0102) |
| **#Medication** | -0.0094 | (-0.0236 , 0.0047) | 0.0033 | (-0.0155 , 0.0221) | 0.0099 | (-0.0113 , 0.0311) | -0.0134 | (-0.0324 , 0.0056) |
| **SumD** | -0.0348 | (-0.0459 , -0.0237)* | -0.0395 | (-0.0510 , -0.0281)* | -0.0339 | (-0.0449 , -0.0229)* | -0.0332 | (-0.0442 , -0.0222)* |
| **SumE** | 0.0002 | (-0.0060 , 0.0064) | -0.0098 | (-0.0147 , -0.0049)* | -0.0096 | (-0.0153 , -0.0039)* | -0.0070 | (-0.0146 , 0.0007) |
| **SumS** | -0.0005 | (-0.0042 , 0.0032) | -0.0040 | (-0.0073 , -0.0007)* | -0.0025 | (-0.0069 , 0.0020) | 0.0014 | (-0.0029 , 0.0058) |
| **Factor3 ~ covariates** | | | | | | | | |
| **Age** | -0.0087 | (-0.0159 , -0.0015)* | -0.0057 | (-0.0285 , 0.0171) | -0.0085 | (-0.0207 , 0.0037) | -0.0114 | (-0.0376 , 0.0147) |
| **Gender** | 0.0044 | (-0.0818 , 0.0906) | 0.3727 | (0.1209 , 0.6246)* | 0.0267 | (-0.1218 , 0.1751) | -0.0343 | (-0.3143 , 0.2457) |
| **Education** | -0.0515 | (-0.1590 , 0.0560) | 0.3347 | (-0.0133 , 0.6827) | 0.0898 | (-0.0768 , 0.2564) | 0.2939 | (-0.0520 , 0.6397) |
| **SumO** | 0.0073 | (0.0007 , 0.0139)* | 0.0084 | (-0.0146 , 0.0315) | 0.0042 | (-0.0089 , 0.0172) | 0.0159 | (-0.0111 , 0.0430) |
| **SumC** | 0.0080 | (-0.0171 , 0.0330) | 0.0147 | (-0.0612 , 0.0905) | -0.0101 | (-0.0365 , 0.0163) | 0.0635 | (0.0014 , 0.1256)* |
| **#Medication** | -0.0191 | (-0.0355 , -0.0027)* | -0.0935 | (-0.1591 , -0.0278)* | -0.1162 | (-0.1618 , -0.0706)* | -0.2878 | (-0.4145 , -0.1611)* |
| **SumD** | -0.0148 | (-0.0276 , -0.0021)* | -0.0552 | (-0.0967 , -0.0137)* | -0.0324 | (-0.0587 , -0.0062)* | -0.0901 | (-0.1465 , -0.0336)* |
| **SumE** | -0.0375 | (-0.0449 , -0.0301)* | -0.1304 | (-0.1480 , -0.1128)* | -0.0669 | (-0.0783 , -0.0554)* | -0.0855 | (-0.1077 , -0.0632)* |
| **SumS** | -0.0068 | (-0.0110 , -0.0026)* | -0.0042 | (-0.0164 , 0.0079) | 0.0022 | (-0.0060 , 0.0105) | -0.0198 | (-0.0360 , -0.0037)* |

* Results are statistically significant with p-value< 0.05

R.C. present regression coefficients ; +:95%C.I.for R.C.

Estimated Correlation Parameters: alpha (0.7938±0.0251) for Factor1

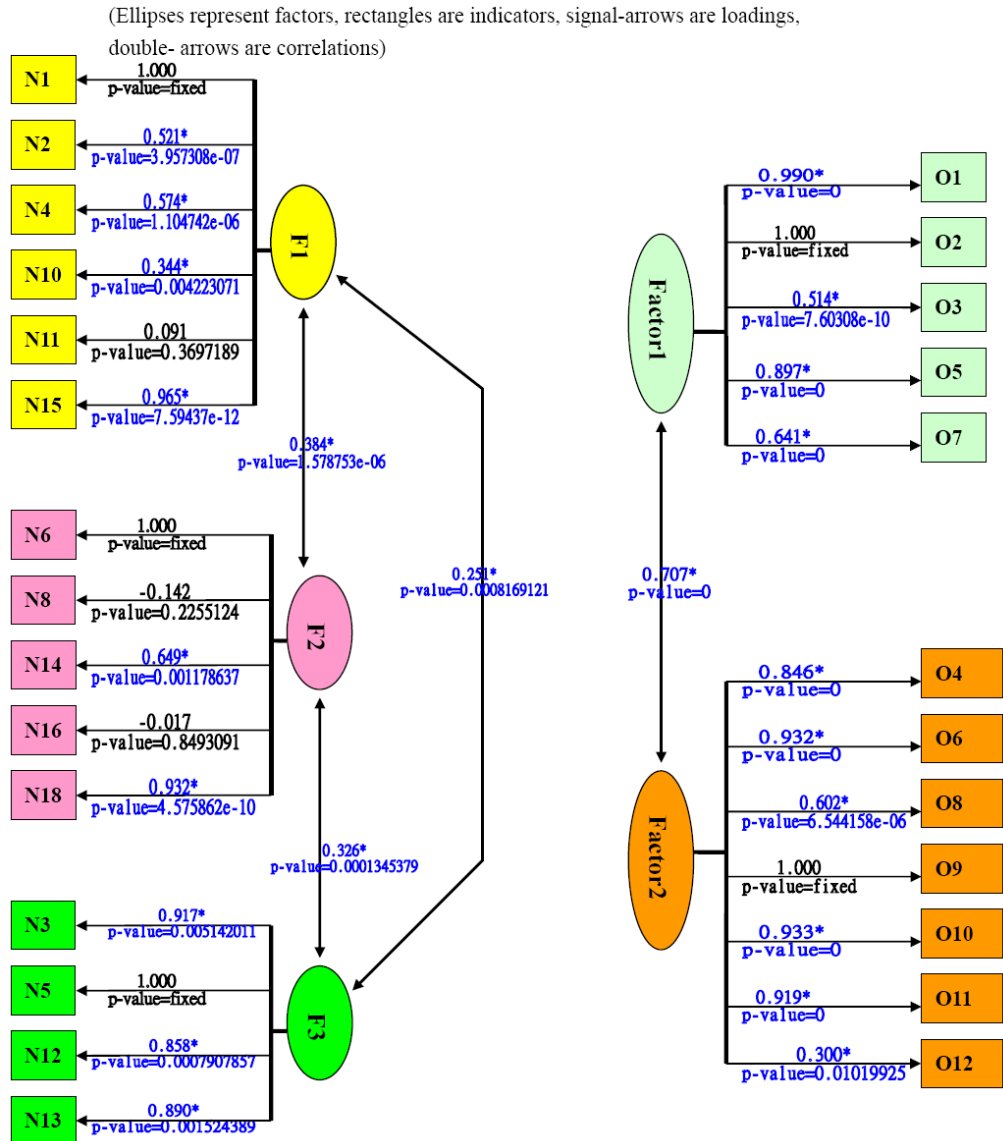; alpha (0.2276±0.0313) for Factor2

; alpha (0.0830±0.0289) for Factor3

Table A.9: The relationship of the difference in factor scores between two adjacent time points with risk factors

| | Time1 ~ Time2 | | Time2 ~ Time3 | | Time3 ~ Time4 | |
| --- | --- | --- | --- | --- | --- | --- |
| | R.C. | 95%C.I.+ | R.C. | 95%C.I.+ | R.C. | 95%C.I.+ |
| **Factor1 score change ~ covariates** | | | | | | |
| Age | 0.0030 | (-0.0041 , 0.0101) | -0.0036 | (-0.0128 , 0.0056) | -0.0011 | (-0.0099 , 0.0077) |
| Gender | -0.0113 | (-0.0926 , 0.0700) | 0.0548 | (-0.0645 , 0.1740) | 0.0804 | (-0.0318 , 0.1930) |
| Education | -0.1306 | (-0.2619 , 0.0006) | 0.0373 | (-0.1300 , 0.2050) | 0.1316 | (-0.0072 , 0.2700) |
| SumO | 0.0027 | (-0.0040 , 0.0095) | 0.0020 | (-0.0069 , 0.0109) | 0.0106 | (-0.0017 , 0.0230) |
| SumC | -0.0170 | (-0.0399 , 0.0060) | 0.0011 | (-0.0253 , 0.0275) | -0.0002 | (-0.0255 , 0.0251) |
| #Medication | 0.0201 | ( 0.0003 , 0.0399)* | -0.0222 | (-0.0662 , 0.0218) | -0.0281 | (-0.0704 , 0.0142) |
| SumD | 0.0026 | (-0.0118 , 0.0169) | 0.0006 | (-0.0198 , 0.0211) | -0.0035 | (-0.0259 , 0.0189) |
| SumE | -0.0084 | (-0.0168 , -0.0001)* | 0.0013 | (-0.0101 , 0.0126) | 0.0019 | (-0.0099 , 0.0138) |
| SumS | 0.0003 | (-0.0044 , 0.0050) | -0.0013 | (-0.0065 , 0.0040) | 0.0017 | (-0.0038 , 0.0073) |
| **Factor2 score change ~ covariates** | | | | | | |
| Age | 0.0006 | (-0.0069 , 0.0082) | 0.0087 | ( 0.0008 , 0.0167)* | -0.0033 | (-0.0104 , 0.0038) |
| Gender | 0.0655 | (-0.0177 , 0.1486) | -0.0771 | (-0.1591 , 0.0048) | -0.0183 | (-0.1028 , 0.0663) |
| Education | -0.0721 | (-0.1783 , 0.0340) | -0.0378 | (-0.1441 , 0.0685) | 0.0472 | (-0.0578 , 0.1522) |
| SumO | 0.0027 | (-0.0041 , 0.0095) | 0.0005 | (-0.0074 , 0.0085) | -0.0017 | (-0.0110 , 0.0075) |
| SumC | -0.0373 | (-0.0651 , -0.0095)* | 0.0209 | (-0.0039 , 0.0456) | -0.0034 | (-0.0201 , 0.0134) |
| #Medication | 0.0127 | (-0.0083 , 0.0338) | 0.0066 | (-0.0212 , 0.0344) | -0.0233 | (-0.0491 , 0.0025) |
| SumD | -0.0047 | (-0.0188 , 0.0095) | 0.0056 | (-0.0096 , 0.0207) | 0.0008 | (-0.0141 , 0.0156) |
| SumE | -0.0100 | (-0.0179 , -0.0021)* | 0.0002 | (-0.0067 , 0.0072) | 0.0026 | (-0.0060 , 0.0112) |
| SumS | -0.0035 | (-0.0076 , 0.0006) | 0.0015 | (-0.0035 , 0.0065) | 0.0039 | (-0.0018 , 0.0096) |
| **Factor3 score change ~ covariates** | | | | | | |
| Age | 0.0030 | (-0.0191 , 0.0252) | -0.0028 | (-0.0272 , 0.0216) | -0.0029 | (-0.0306 , 0.0248) |
| Gender | 0.3683 | ( 0.1161 , 0.6205)* | -0.3461 | (-0.6430 , -0.0492)* | -0.0610 | (-0.3442 , 0.2223) |
| Education | 0.3862 | ( 0.0428 , 0.7297)* | -0.2450 | (-0.6145 , 0.1246) | 0.2041 | (-0.1588 , 0.5670) |
| SumO | 0.0012 | (-0.0220 , 0.0243) | -0.0043 | (-0.0320 , 0.0234) | 0.0118 | (-0.0172 , 0.0408) |
| SumC | 0.0067 | (-0.0709 , 0.0843) | -0.0247 | (-0.1030 , 0.0535) | 0.0736 | ( 0.0113 , 0.1359)* |
| #Medication | -0.0743 | (-0.1407 , -0.0079)* | -0.0228 | (-0.1012 , 0.0557) | -0.1716 | (-0.3007 , -0.0425)* |
| SumD | -0.0404 | (-0.0832 , 0.0024) | 0.0228 | (-0.0239 , 0.0694) | -0.0576 | (-0.1216 , 0.0064) |
| SumE | -0.0929 | (-0.1117 , -0.0742)* | 0.0636 | ( 0.0428 , 0.0844)* | -0.0186 | (-0.0437 , 0.0064) |
| SumS | 0.0026 | (-0.0099 , 0.0151) | 0.0065 | (-0.0085 , 0.0215) | -0.0221 | (-0.0400 , -0.0042)* |

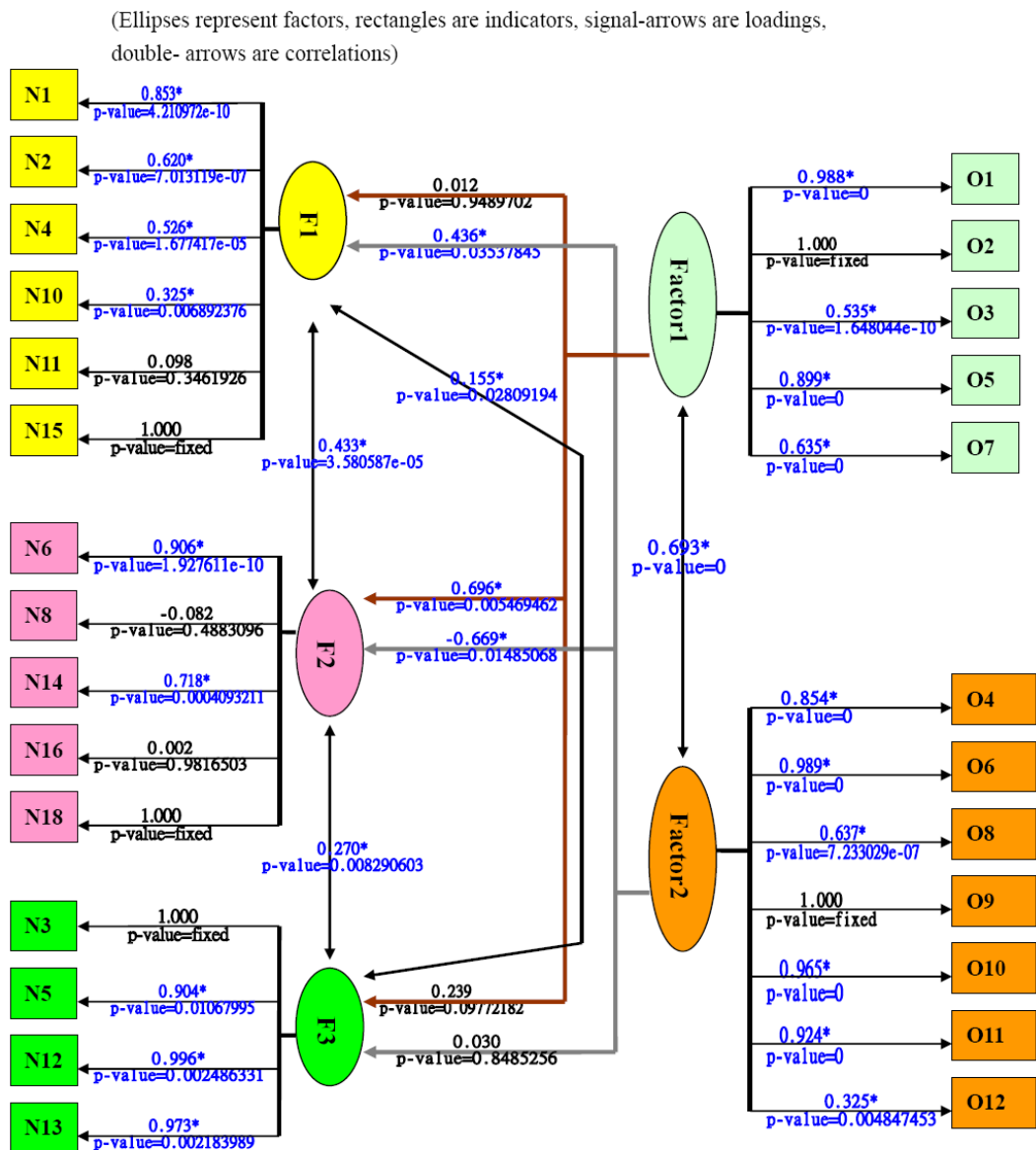* Results are statistically significant with p-value<0.05

R.C. present regression coefficients ; +:95%C.I.for R.C.

Figure A.1: The Measurement models for oral health and nutrition

(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings,
double- arrows are correlations)

| | |
|---|---|
| **N1** | 1.000 / p-value=fixed |
| **N2** | 0.521* / p-value=3.957308e-07 |
| **N4** | 0.574* / p-value=1.104742e-06 |
| **N10** | 0.344* / p-value=0.004223071 |
| **N11** | 0.091 / p-value=0.3697189 |
| **N15** | 0.965* / p-value=7.59437e-12 |

**F1**

0.384* / p-value=1.578753e-06

0.251* / p-value=0.0008169121

| | |
|---|---|
| **N6** | 1.000 / p-value=fixed |
| **N8** | -0.142 / p-value=0.2255124 |
| **N14** | 0.649* / p-value=0.001178637 |
| **N16** | -0.017 / p-value=0.8493091 |
| **N18** | 0.932* / p-value=4.575862e-10 |

**F2**

0.326* / p-value=0.0001345379

| | |
|---|---|
| **N3** | 0.917* / p-value=0.005142011 |
| **N5** | 1.000 / p-value=fixed |
| **N12** | 0.858* / p-value=0.0007907857 |
| **N13** | 0.890* / p-value=0.001524389 |

**F3**

**Factor1**

| | |
|---|---|
| 0.990* / p-value=0 | **O1** |
| 1.000 / p-value=fixed | **O2** |
| 0.514* / p-value=7.60308e-10 | **O3** |
| 0.897* / p-value=0 | **O5** |
| 0.641* / p-value=0 | **O7** |

0.707* / p-value=0

**Factor2**

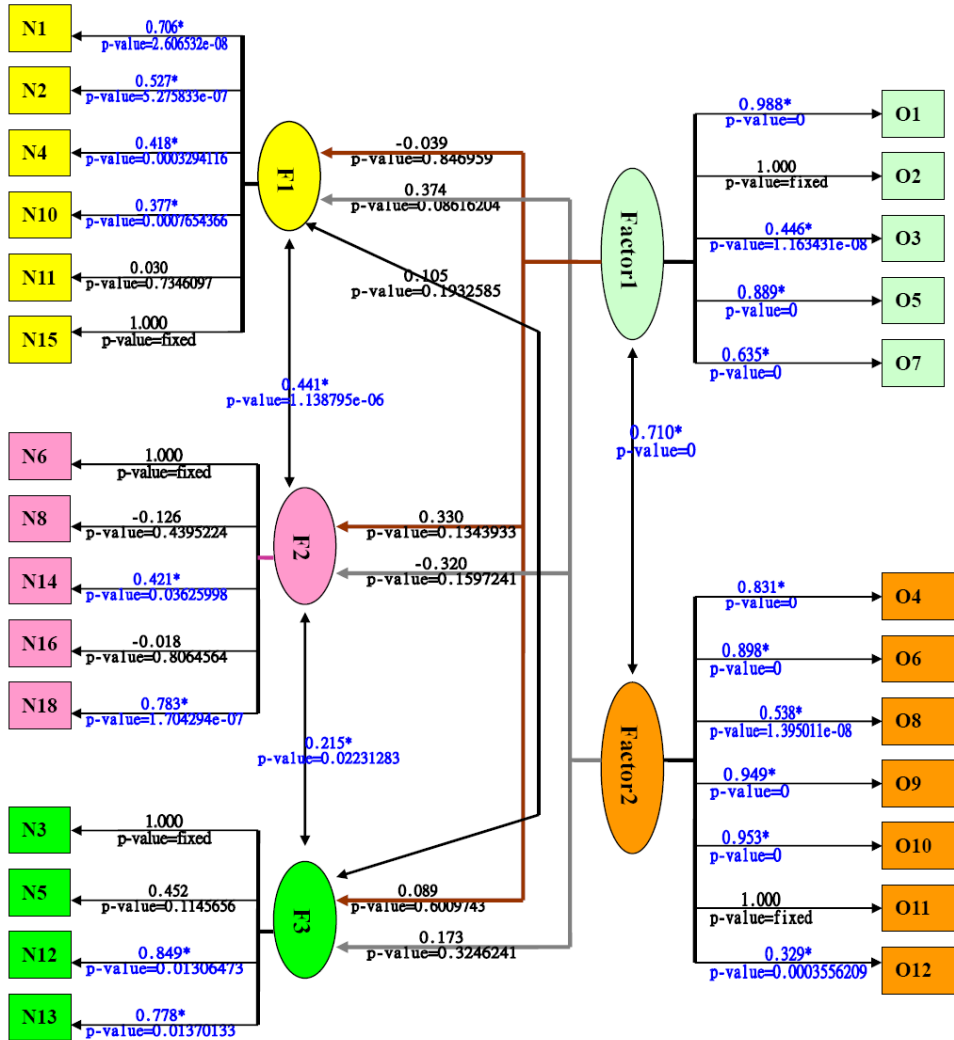| | |
|---|---|
| 0.846* / p-value=0 | **O4** |
| 0.932* / p-value=0 | **O6** |
| 0.602* / p-value=6.544158e-06 | **O8** |
| 1.000 / p-value=fixed | **O9** |
| 0.933* / p-value=0 | **O10** |
| 0.919* / p-value=0 | **O11** |
| 0.300* / p-value=0.01019925 | **O12** |

* Results are statistically significant with p-value < 0.05

Figure A.2: Structural equation modeling: Combining measurement and structural models



(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings, double- arrows are correlations)

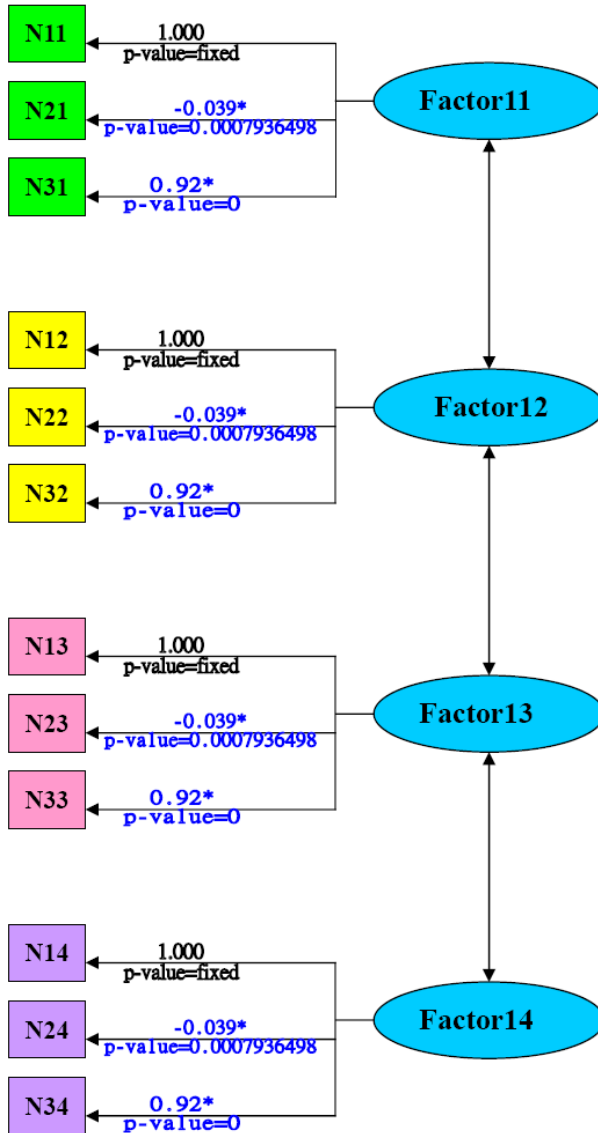* Results are statistically significant with p-value < 0.05

Figure A.3: Structural equation modeling: Combining measurement and structural models. The relationship is adjusted for risk factors (Age, Gender, Education, SumC, #Medication, SumD, SumE, SumS)



(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings, double- arrows are correlations)

* Results are statistically significant with p-value < 0.05

Figure A.4: CFA for Factor1 at different time points

(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings, double- arrows are correlations)

Factorij: i presents factor and $i = 1,2,3$ ; j presents time point and $j = 1,2,3,4$

Nij: i presents indicator and $i = 1,2,...,18$ ; j present time point and $j = 1,2,3,4$

| N11 | | 1.000 |
| N21 | | p-value=fixed |
| N31 | | -0.039* p-value=0.0007936498 |
| | | 0.92* p-value=0 |

Factor11

| N12 | | 1.000 |
| | | p-value=fixed |
| N22 | | -0.039* p-value=0.0007936498 |
| N32 | | 0.92* p-value=0 |

Factor12

| N13 | | 1.000 |
| | | p-value=fixed |
| N23 | | -0.039* p-value=0.0007936498 |
| N33 | | 0.92* p-value=0 |

Factor13

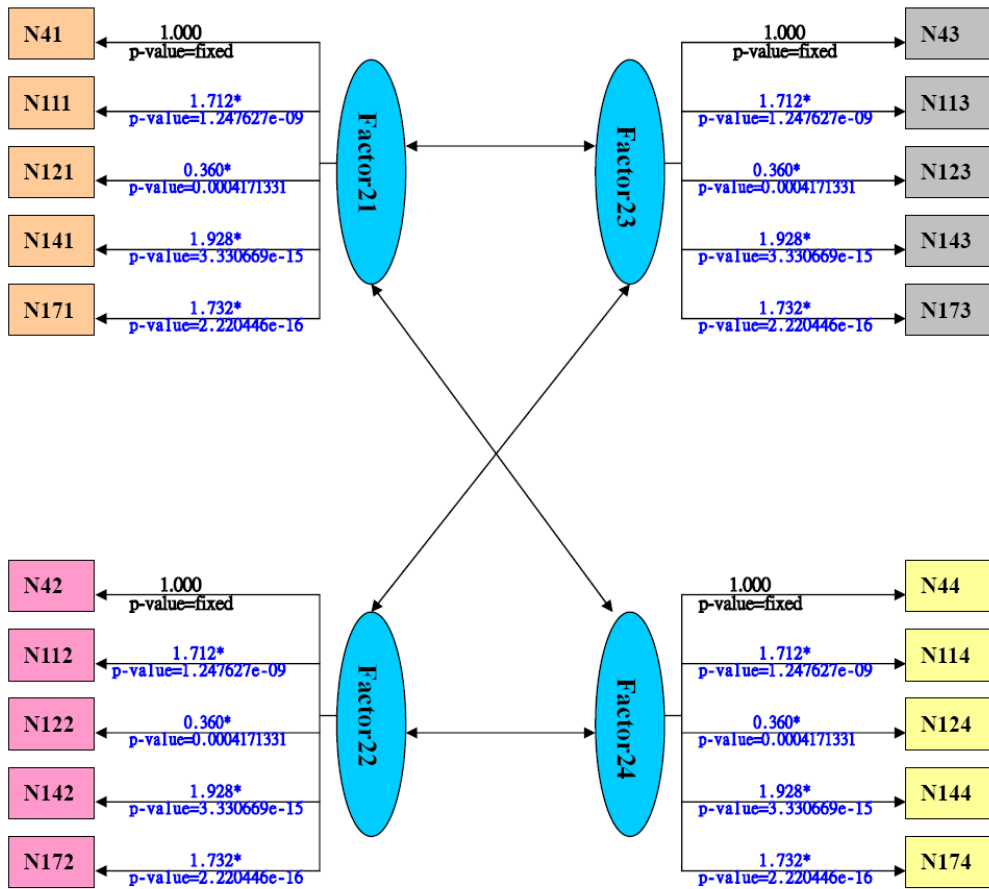| N14 | | 1.000 |
| | | p-value=fixed |
| N24 | | -0.039* p-value=0.0007936498 |
| N34 | | 0.92* p-value=0 |

Factor14

* Results are statistically significant with p-value < 0.05

## Figure A.5: CFA for Factor2 at different time points

(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings,
double- arrows are correlations)
Factorij: i presents factor and $i = 1,2,3$ ; j presents time point and $j = 1,2,3,4$
Nij: i presents indicator and $i = 1,2,....,18$ ; j present time point and $j = 1,2,3,4$
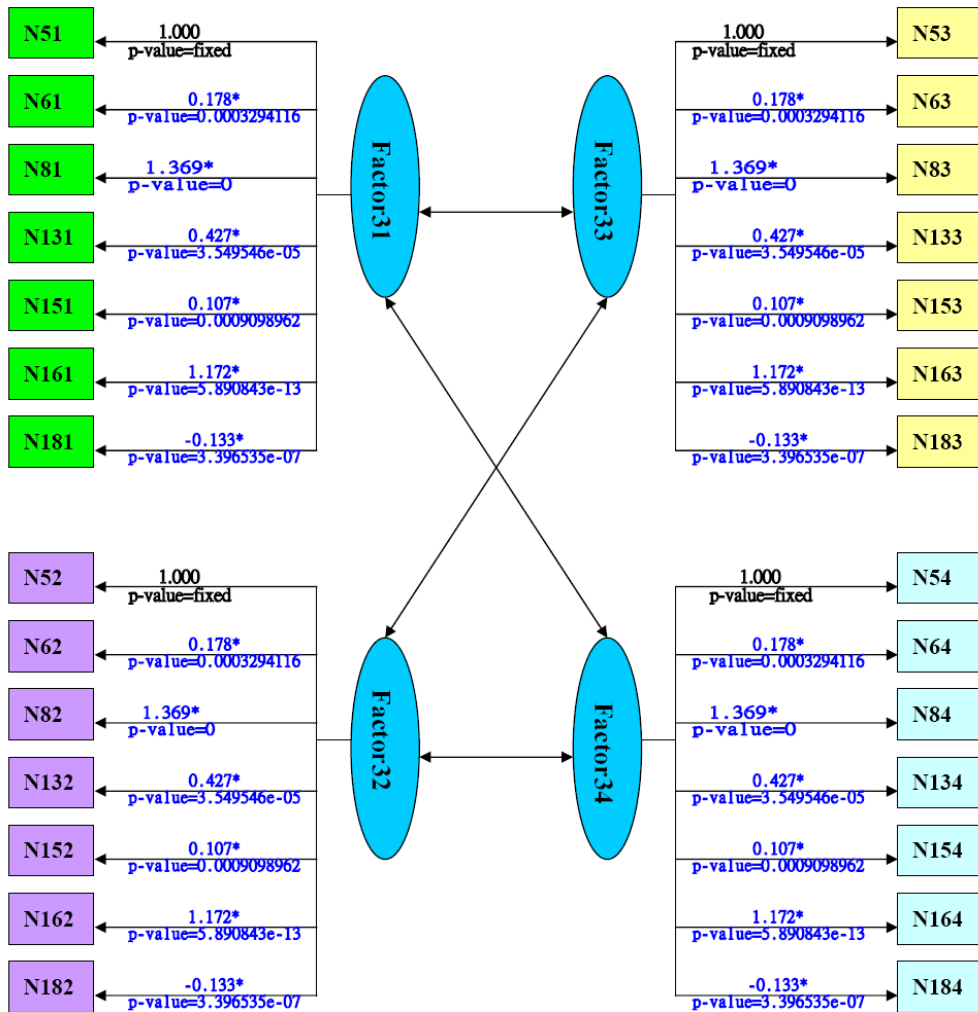


* Results are statistically significant with p-value < 0.05

# Figure A.6: CFA for Factor3 at different time points

(Ellipses represent factors, rectangles are indicators, signal-arrows are loadings, double- arrows are correlations)

Factorij: i presents factor and $i = 1,2,3$ ; j presents time point and $j = 1,2,3,4$

Nij: i presents indicator and $i = 1,2,...,18$ ; j present time point and $j = 1,2,3,4$



* Results are statistically significant with p-value < 0.05