

Analysis of Longitudinal Survey Data

Introduction to Generalized Estimating Equations
with Examples from the ITC Survey

Pete Driezen

June 13, 2016

Introduction

- To date, an ITC Survey has been conducted in 23 countries around the world, countries inhabited by
 - 50% of the world's population
 - 60% of the world's smokers and
 - 70% of the world's tobacco users
- ITC surveys:
 - have an international scope
 - use quasi-experimental designs
 - employ representative samples of smokers
 - measure outcomes for multiple tobacco control policies
 - measure intermediate outcomes (i.e., what a policy should be changing before an individual changes behaviour)
 - include psychosocial mediators and moderators that shed light on how and why the policy works

Features of the ITC Surveys

- Probability-based sampling designs
 - nationally or regionally representative samples
 - stratified designs (e.g., ITC 4 Country, ITC Netherlands)
 - multi-stage cluster designs (e.g., ITC Bangladesh, ITC Zambia)
- Interviewing methods appropriate for the design
 - telephone and/or web (high income countries)
 - face-to-face (LMICs)
- Longitudinal design
 - cohort members followed over time
 - respondents lost to attrition are replenished

Advantages of the ITC Design

- ITC research design takes advantage of natural experiments
 - pre-post designs: policy change within single countries
 - pre-post comparison designs: policy change in one country but not in others (e.g., 4 Country Survey, ITC Southeast Asia)
- ITC permits examination of:
 - change within individuals over time
 - mediation effects
- Estimation of population level temporal trends
 - multiple surveys conducted within countries over several years
- Cross-country comparisons
 - an ITC survey has been conducted in 23 countries around the world

Analytic Challenges

- Generally, observations from an ITC survey are **not** independent and identically distributed (IID):
 1. **Multistage cluster designs:** respondents sampled from the same areas are more alike than respondents sampled from different areas
 2. **Repeated measures:** respondents surveyed at multiple points in time (within-subject correlation)
 3. **Time-in-sample effects:** mix of re-contact and replenishment respondents (respondents surveyed more than once often differ systematically from those surveyed for the first time)
- These features need to be accounted for in the analysis of ITC data to:
 - produce correct variance estimates (and test statistics)
 - produce unbiased population-level estimates (in the case of time-in-sample effects)

Generalized Estimating Equations

- GEE models take the form of regression models with correlation within subjects
- Can be estimated for different types of outcomes using a “link” function:

Response	Distribution	Link
Continuous	Normal	identity
Binary	Binomial	logit
Nominal	Multinomial	generalized logit
Count	Poisson	log

- GEE models are “marginal” or population-averaged models
- Practically, this means that GEE parameter estimates refer to an **average** person in the population
- In a policy context, making inferences about what happens on average in the population is acceptable, so GEE methods are useful for ITC data

Accounting for Correlated Responses

- Correlated data arise through:
 - multi-stage cluster sampling designs
 - nesting of observations within larger units (similar to cluster sampling), e.g., students nested within schools
 - repeated measures (a type of “nesting”), i.e., multiple measurements on the same individual at different time points
- Statistical methods need to account for this correlation (see [Hanley et al., 2003](#))
- GEE models allow specification of different types of correlation structures. The type used depends on the nature of the data (e.g., longitudinal data vs. clustered data)
- However, GEE models are robust to misspecification of the correlation structure (i.e., parameter estimates are not dramatically affected by specifying the incorrect correlation)

General Analytic Approach

- Estimating a GEE model is an iterative process that starts with maximum-likelihood estimation of regression parameters (β 's)
- Variances [$\text{var}(\beta)$] estimated using an appropriate link function
- $\text{var}(\beta)$ are multiplied against a working matrix of correlation coefficients that corrects for the correlation within subjects
- Nature of the data helps determine correlation structure
- Procedure repeats until the change in parameter estimates from one iteration to the next approaches zero
- Using the correct form of correlation increases estimation efficiency, but estimates of regression parameters are consistent even if correlation structure is incorrect
- Working correlation is usually estimated from the data
- More complex correlation structures mean more parameters are estimated, which may lead to estimation problems in some cases

Types of Correlation Structure

- **Independent:**
 - assumes observations are not correlated
- **Exchangeable (compound symmetry):**
 - assumes within-subject observations are equally correlated
 - appropriate when observations are nested within units, e.g., students in schools or puppies in litters
- **Auto-regressive:**
 - correlation is an exponential function of the lag period
 - e.g., stronger correlation between within-subject responses measured at times 1 & 2 than at times 1 & 3
- **Unstructured:**
 - free estimation of the correlation (most complex structure)
 - estimates all possible correlations
- Unstructured and exchangeable working correlation structures most commonly used with ITC data

Common Working Correlation Models

- See Horton & Lipsitz for general forms (p 161)
- Concretely, for data consisting of measurements made on the same respondents for each of three time points, then:

Unstructured correlation matrix:

$$\begin{matrix} 1 & \rho_{1,2} & \rho_{1,3} \\ \rho_{1,2} & 1 & \rho_{2,3} \\ \rho_{1,3} & \rho_{2,3} & 1 \end{matrix}$$

Exchangeable correlation matrix:

$$\begin{matrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{matrix}$$

Estimating GEE Models with ITC Data

- With data arising from ITC surveys, often need to account for multiple issues arising from these complex survey data
 - sampling weights
 - sampling design
 - repeated measures
 - time-in-sample
- Observations for analysis — choice affects which sampling weight is used for model estimation:
 - respondents present in **all** waves: longitudinal weight
 - respondents present in **any** wave: cross-sectional weight
- In either case, the sampling weight **must** remain constant within respondents for all time points included in the analysis
- Explanatory variables may be time-invariant (fixed) or time-varying. Socio-demographic covariates are typically treated as time-invariant.

Software to Estimate GEE Models

- Most commonly use SAS and SAS-callable SUDAAN for estimation of GEE models using ITC data
- SAS:
 - “proc genmod” allows estimation of linear and logistic GEE models
- SUDAAN: several procedures
 - running a GEE model depends on which options are passed to the procedure
 - “proc regress” (linear models), “proc rlogist” (logistic models), “proc multilog” (ordinal or multinomial models)
- Other software:
 - Stata (“xtgee”)
 - R (requires additional packages: “gee”, “geepack” or “multgee”)
- SUDAAN is the only package that can account for **both** the sampling design and repeated measures

Data Arrangement

- Estimation of GEE models requires data to be arranged in a “long” format rather than the usual “wide” format
- Wide:
 - each row represents a single, unique respondent
 - multiple measurements for the same variable are stored in separate columns (e.g., “NoticeHWL1”, “NoticeHWL2”, ..., “NoticeHWL7”, etc.)
- Long:
 - each respondent has multiple rows of data, one for each time point
 - respondents are indexed by a “time” or “wave” indicator
 - sampling weight is constant within unique individuals

“Long” Data Arrangement

long

Filter

	uniqid	country	cohort	wave	tis	strata	xwt	sex	ageGroup	noticeHWL
1	10137400042	1	1	1	1	101	0.6657718	2	1	0
2	10153000050	1	1	1	1	101	0.6843693	1	2	1
3	10153000050	1	1	2	2	101	0.6843693	1	2	1
4	10153000050	1	1	3	3	101	0.6843693	1	2	1
5	10153000050	1	1	4	4	101	0.6843693	1	2	1
6	10153000050	1	1	5	5	101	0.6843693	1	2	1
7	10153000050	1	1	6	6	101	0.6843693	1	2	1
Multiple rows of data for unique respondents										
9	10172700060	1	1	2	2	101	1.2258966	1	1	0
10	10181400064	1		1	1	101	0.7771534		3	0
11	10181400064	1		2	2	101	0.7771534		3	0
12	10181400064	1	1	3	3	101	0.7771534	2	3	0
13	10181400064	1	1	4	4	101	0.7771534	2	3	0
14	10181400064	1	1	5	5	101	0.7771534	2	3	1
15	10181400064	1	1	6	6	101	0.7771534	2	3	0
16	10193400072	1	1	1	1	101	0.6843693	1	2	1
17	10193400072	1	1	2	2	101	0.6843693	1	2	0
18	10193400072	1	1	3	3	101	0.6843693	1	2	0
19	10193400072	1	1	4	4	101	0.6843693	1	2	0
20	10193400072	1	1	5	5	101	0.6843693	1	2	0
21	10193400072	1	1	8	8	101	0.6843693	1	2	0
22	10235700085	1	1	1	1	101	0.6129483	1	1	0

Showing 1 to 23 of 28,080 entries

Data Preparation

- Preparing for GEE analysis in SAS involves a bit of data management
- Essentially requires ensuring multiple waves of data contain identically named variables and an indicator variable for wave (or “time”). Data can then be “stacked”
- Somewhat easier in Stata (using the “reshape” command) and R (using “melt” and “cast” functions in the “reshape2” package)

Package	Topic
SAS	Data management
	Stacking data
Stata	Data management
	Reshape: long to wide
	Reshape: Wide to long
R	Starter kit
	“reshape2” package

Example: Thailand Warning Labels

Background

- Pictorial health warning labels (HWLs) were introduced in Thailand on March 25, 2006 (text only labels from 1997 to that time)
- Label size was increased from 33% to 50% of the front & back of the pack
- Data from the ITC Southeast Asia Survey were used to evaluate the effect of the pictorial HWLs on (a) salience of health warnings, (b) cognitive reactions to health warnings and (c) behavioural reactions health warnings

Methods

- Data from the first 3 waves of the ITC Southeast Asia Survey
- Malaysia used as comparison group (HWLs were unchanged)
- 6,287 unique smokers participating in one or more waves (Malaysia n = 3,220; Thailand n = 3,067)
- New HWLs introduced just after Wave 1 was completed

Thailand Warning Labels

Outcome measures

- Salience: noticing and closely reading HWLs “Often” or “Very often”
- Cognitive reactions: thinking about the harms of smoking and quitting “a lot”
- Behavioural reactions: forgoing a cigarette “at least once”

Covariates

- Standard socio-demographic measures, including sex, age group, residence (urban vs. rural), income, education, ethnicity and wave of recruitment (similar to time-in-sample)
- Behavioural measures: smoking status (daily vs. non-daily), cigarettes smoked per day, use of RYO tobacco

Thailand Warning Labels

Analysis

- Estimate temporal changes in salience of and reactions to HWL in Thailand and Malaysia
 - Were changes larger in Thailand following introduction of new HWL compared to Malaysia where HWL remained the same?
- Test whether salience of and reactions to HWL increased significantly from Wave 1 to Wave 2 in Thailand (and whether changes were maintained in Wave 3)
- Test whether trends differed between countries
- Re-analysis of Yong et al., 2013, using GEE methods in both SAS and SUDAAN to estimate trends and test effects
- Key to testing differences over time between countries is the incorporation of a country X wave interaction term in the GEE model

Estimating a GEE Model using SAS

- Weighted GEE models can be estimated using SAS, but it is not possible to account for the complex sampling design
 - more problematic for multi-stage cluster designs, where neighbourhoods or villages form the primary sampling units
 - for stratified designs, such as ITC 4 Country or ITC Netherlands, estimating GEE models in SAS (ignoring the strata) is reasonable
 - approach used in evaluation of the Ireland smoke-free bars policy (Fong et al., 2006)
- GEE models estimated using SUDAAN account for both the complex sampling design **and** repeated measures
 - however, only have a choice of two correlation structures: independent or exchangeable
 - since GEE models are robust to misspecification of the correlation structure, estimates from SUDAAN are generally reasonable

Unadjusted GEE Model

```
/* Country-specific estimates by wave. Note: cohort = wave of recruitment */
/* SAS */
proc genmod data = hwl order = internal desc;
  class unqid sex agegrp urban income educ ethnic daily cpd ryob
    cohort country wave / param = glm;
  weight xwt;
  model wlnotice = country wave country*wave /
    dist = bin link = logit type3 wald;
  lsmeans country*wave / om ilink cl;
  repeated subject = unqid / corrw type = exch withinsubject = wave;
run;

/* SUDAAN */
proc rlogist data = hwl r = exchangeable semethod = zeger;
  setenv decwidth = 4;
  nest strata psu unqid / psulev = 2;
  weight xwt;
  class sex agegrp urban income educ ethnic daily cpd ryob
    cohort country wave / nofreq dir = ascending;
  model wlnotice = country wave country*wave;
  predmarg country*wave;
run;
```

```

proc genmod data = hwl order = internal desc;
  class unqid sex agegrp urban income educ ethnic daily cpd ryob
    cohort country wave / param = glm;
  weight xwt;
  model wlnotice = country wave country*wave /
    dist = bin link = logit type3 wald;
  lsmeans country*wave / om ilink cl;
  repeated subject = unqid / corrw type = exch withinsubject = wave;
run;

```

- PROC GENMOD can handle a variety of generalized linear models
- “wlnotice” → binary variable. To estimate a logistic regression model, need to use the “dist = bin link = logit” option. This specifies that the outcome follows a binomial distribution. A logit link function is used to estimate the model
- “param = glm” → requests “glm” coding of categorical variables specified on the “class” statement (needed for “lsmeans”)
- “lsmeans” → “least squares means:”
 - marginal estimates or group means after controlling for other covariates

```

proc genmod data = hwl order = internal desc;
  class unqid sex agegrp urban income educ ethnic daily cpd ryob
    cohort country wave / param = glm;
  weight xwt;
  model wlnotice = country wave country*wave /
    dist = bin link = logit type3 wald;
  lsmeans country*wave / om ilink cl;
  repeated subject = unqid / corrw type = exch withinsubject = wave;
run;

```

- “lsmeans” →
 - by default, SAS estimates these marginal effects using a hypothetical balanced population (not what we want). Instead, need adjusted estimates that reflect the “global average respondent” or a respondent who possesses the average values of all covariates used in model estimation
 - ∴ use the “om” or “obsmargin” (“observed margins”)
 - “ilink” → inverse link to estimate probabilities (instead of log odds ratios)
 - “cl” → requests confidence intervals

```

proc genmod data = hwl order = internal desc;
  class unqid sex agegrp urban income educ ethnic daily cpd ryob
    cohort country wave / param = glm;
  weight xwt;
  model wlnotice = country wave country*wave /
    dist = bin link = logit type3 wald;
  lsmeans country*wave / om ilink cl;
  repeated subject = unqid / corrw type = exch withinsubject = wave;
run;

```

- “repeated subject = unqid” → specifies how observations are identified as repeated
 - “corrw” → prints the estimated working correlation matrix
 - “type = exch” → exchangeable working correlation. Other options:
 - “unstr” = unstructured
 - “ar” = first-order autoregressive
 - “ind” = independent
 - “fixed” = user-specified matrix
 - “withinsubject = wave” → defines the order of observations within subjects

```

proc rlogist data = hwl r = exchangeable semethod = zeger;
  setenv decwidth = 4;
  nest strata psu uniqid / psulev = 2;
  weight xwt;
  class sex agegrp urban income educ ethnic daily cpd ryob
        cohort country wave / nofreq dir = ascending;
  model wlnotice = country wave country*wave;
  predmarg country*wave;
run;

```

- Most regression procedures in SUDAAN can estimate a GEE model with complex survey data → requires “r = exchangeable semethod = zeger” options on the model statement
- “proc rlogistic” estimates a logistic regression model
- “nest strata psu uniqid / psulev = 2” specifies the sampling design
 - first variable specifies the sampling strata
 - second variable specifies the primary sampling unit
 - third variable represents the repeated measures, but this requires correct identification of PSUs using the “psulev = 2” option (i.e., second variable represents PSUs)
- “class” defines categorical variables
- “model” statement → as in SAS


```
proc rlogist data = hwl r = exchangeable semethod = zeger;  
  setenv decwidth = 4;  
  nest strata uniqid;  
  weight xwt;  
  class sex agegrp urban income educ ethnic daily cpd ryob  
    cohort country wave / nofreq dir = ascending;  
  model wlnotice = country wave country*wave;  
  predmarg country*wave;  
run;
```

- “predmarg” produces “predicted marginals” → similar to “lsmeans” in SAS but more appropriate (“condmarg” or “conditional marginals” produces the same estimates as “lsmeans” in SAS)
- For a logistic model, “predicted marginals” are on the probability scale (proc genmod in SAS requires specification of the “inverse link function”)

----- S A S -----

Exchangeable Working Correlation: 0.1894

country*wave Least Squares Means

country	wave	Mean	Standard Error of Mean	Lower Mean	Upper Mean
1. Malaysia	1	0.5722	0.01704	0.5385	0.6052
1. Malaysia	2	0.5076	0.02005	0.4684	0.5468
1. Malaysia	3	0.5426	0.01783	0.5075	0.5773
2. Thailand	1	0.6229	0.01243	0.5982	0.6469
2. Thailand	2	0.7001	0.01198	0.6761	0.7231
2. Thailand	3	0.7389	0.01097	0.7168	0.7598

----- S U D A A N -----

Working Correlations: Exchangeable (rho = 0.1879)

Predicted Marginal #1	Predicted Marginal	SE	Lower 95% Limit	Upper 95% Limit

country, WAVE				
1. Malaysia, 1	0.5722	0.0402	0.4902	0.6504
1. Malaysia, 2	0.5076	0.0349	0.4378	0.5771
1. Malaysia, 3	0.5425	0.0347	0.4725	0.6109
2. Thailand, 1	0.6229	0.0238	0.5739	0.6694
2. Thailand, 2	0.7001	0.0218	0.6547	0.7420
2. Thailand, 3	0.7389	0.0188	0.6994	0.7748

Multivariable Model (Adjusted Estimates)

```
proc rlogist data = hwl r = exchangeable semethod = zeger;
  nest strata psu unqid / psulev = 2;
  weight xwt;
  class sex agegrp urban income educ ethnic daily cpd ryob
        cohort country wave / nofreq dir = ascending;
  model wlnotice = sex agegrp urban income educ ethnic daily
        cpd ryob cohort country wave country*wave;
  predmarg country*wave;
  /* Tests difference in predicted marginals */
  pred_eff country = (1 0) * wave = (-1 1 0) / name = "Malaysia, W2 vs W1";
  pred_eff country = (1 0) * wave = (-1 0 1) / name = "Malaysia, W3 vs W1";
  pred_eff country = (1 0) * wave = (0 -1 1) / name = "Malaysia, W3 vs W2";
  pred_eff country = (0 1) * wave = (-1 1 0) / name = "Thailand, W2 vs W1";
  pred_eff country = (0 1) * wave = (-1 0 1) / name = "Thailand, W3 vs W1";
  pred_eff country = (0 1) * wave = (0 -1 1) / name = "Thailand, W3 vs W2";
  /* Can also estimate odds ratios for the effects */
  effects wave = (-1 1 0) / country = 1 exp name = "Malaysia, W2 vs W1";
  effects wave = (-1 0 1) / country = 1 exp name = "Malaysia, W3 vs W1";
  effects wave = (0 -1 1) / country = 1 exp name = "Malaysia, W3 vs W1";
  effects wave = (-1 1 0) / country = 2 exp name = "Thailand, W2 vs W1";
  effects wave = (-1 0 1) / country = 2 exp name = "Thailand, W3 vs W1";
  effects wave = (0 -1 1) / country = 2 exp name = "Thailand, W3 vs W1";
run;
```

Predicted Marginal (Country, Wave)	Predicted Marginal	SE	Lower 95% Limit	Upper 95% Limit
Unadjusted				
Malaysia, 1	0.5722	0.0402	0.4902	0.6504
Malaysia, 2	0.5076	0.0349	0.4378	0.5771
Malaysia, 3	0.5425	0.0347	0.4725	0.6109
Thailand, 1	0.6229	0.0238	0.5739	0.6694
Thailand, 2	0.7001	0.0218	0.6547	0.7420
Thailand, 3	0.7389	0.0188	0.6994	0.7748
Adjusted				
Malaysia, 1	0.4993	0.0370	0.4256	0.5731
Malaysia, 2	0.4487	0.0306	0.3883	0.5106
Malaysia, 3	0.5043	0.0285	0.4471	0.5614
Thailand, 1	0.6233	0.0243	0.5735	0.6707
Thailand, 2	0.7497	0.0197	0.7081	0.7872
Thailand, 3	0.8033	0.0165	0.7680	0.8343

/* Tests Differences Within Countries Between Waves */

Contrast	PREDMARG	SE	T-Stat	P-value	OR	Lower	Upper
Malaysia, W2 vs W1	-0.0507	0.0467	-1.0861	0.2828	0.7964	0.5230	1.2127
Malaysia, W3 vs W1	0.0050	0.0486	0.1020	0.9192	1.0225	0.6592	1.5860
Malaysia, W3 vs W2	0.0556	0.0410	1.3576	0.1808	1.2839	0.8856	1.8613
Thailand, W2 vs W1	0.1264	0.0243	5.2105	0.0000	1.9573	1.5081	2.5404
Thailand, W3 vs W1	0.1799	0.0267	6.7494	0.0000	2.7671	2.0346	3.7632
Thailand, W3 vs W2	0.0536	0.0191	2.8012	0.0073	1.4137	1.1037	1.8108

Between Country Differences

- Recap:
 - estimated percentage of smokers noticing warning labels “often/very often” (unadjusted and adjusted estimates)
 - estimated change in support within each country from:
 - Wave 1 to Wave 2
 - Wave 2 to Wave 3
 - Wave 1 to Wave 3
 - changes estimated as differences in percentages and odds ratios
 - tested differences
- However, it’s also of interest to test temporal differences between countries
- In other words, is the change in support in Thailand between Waves 1 and 2 different than the change in support in Malaysia during this time period?
- In SUDAAN, these differences are also estimated using the `pred_eff` statement

```

proc rlogist data = hwl r = exchangeable semethod = zeger;
  nest strata psu unqid / psulev = 2;
  weight xwt;
  class sex agegrp urban income educ ethnic daily cpd ryob
        cohort country wave / nofreq dir = ascending;
  model wlnotice = sex agegrp urban income educ ethnic daily
        cpd ryob cohort country wave country*wave;
  /* Predicted marginals (adjusted probabilities) */
  predmarg country*wave;
  /* Tests difference in predicted marginals */
  pred_eff country = (1 0) * wave = (-1 1 0) / name = "Malaysia, W2 vs W1";
  pred_eff country = (1 0) * wave = (-1 0 1) / name = "Malaysia, W3 vs W1";
  pred_eff country = (1 0) * wave = (0 -1 1) / name = "Malaysia, W3 vs W2";
  pred_eff country = (0 1) * wave = (-1 1 0) / name = "Thailand, W2 vs W1";
  pred_eff country = (0 1) * wave = (-1 0 1) / name = "Thailand, W3 vs W1";
  pred_eff country = (0 1) * wave = (0 -1 1) / name = "Thailand, W3 vs W2";
  /* Relative change in support over time between countries */
  pred_eff country = (-1 1) * wave = (-1 1 0) / name = "TH W2-W1 vs MY W2-W1";
  pred_eff country = (-1 1) * wave = (-1 0 1) / name = "TH W3-W1 vs MY W3-W1";
  pred_eff country = (-1 1) * wave = (0 -1 1) / name = "TH W3-W2 vs MY W3-W2";
run;

```

Predicted Marginal #1	Predicted Marginal	SE	Lower 95% Limit	Upper 95% Limit
country, WAVE				
Malaysia, 1	0.4993	0.0370	0.4256	0.5731
Malaysia, 2	0.4487	0.0306	0.3883	0.5106
Malaysia, 3	0.5043	0.0285	0.4471	0.5614
Thailand, 1	0.6233	0.0243	0.5735	0.6707
Thailand, 2	0.7497	0.0197	0.7081	0.7872
Thailand, 3	0.8033	0.0165	0.7680	0.8343

Contrasted Predicted Marginal #1	PREDMARG Contrast	SE	T-Stat	P-value
Malaysia, W2 vs W1	-0.0507	0.0467	-1.0861	0.2828
Malaysia, W3 vs W1	0.0050	0.0486	0.1020	0.9192
Malaysia, W3 vs W2	0.0556	0.0410	1.3576	0.1808
Thailand, W2 vs W1	0.1264	0.0243	5.2105	0.0000
Thailand, W3 vs W1	0.1799	0.0267	6.7494	0.0000
Thailand, W3 vs W2	0.0536	0.0191	2.8012	0.0073
TH W2-W1 vs MY W2-W1	0.1771	0.0526	3.3635	0.0015
TH W3-W1 vs MY W3-W1	0.1750	0.0549	3.1885	0.0025
TH W3-W2 vs MY W3-W2	-0.0021	0.0443	-0.0469	0.9628

Multiple Comparisons?

- GEE model in SUDAAN used to test effects of interest:
 - Does the warning label policy change in Thailand between wave 1 and wave 2 significantly affect warning label salience among smokers?
 - Need to know whether any change is likely attributable to the new warning labels → requires use of a comparison group, in this case, Malaysia
 - Also want to control for important differences between countries (demographic and smoking behaviour covariates)
 - ∴ Requires a “country X wave” interaction effect in the GEE model (overall interaction is significant, Wald $F_{2,49} = 8.68$, $p < 0.001$)
- Given overall interaction was significant, explored several effects of interest:
 - 3 within country differences (6 tests) & 3 between country differences = 9 separate statistical tests
 - Bonferroni or False Discovery Rate adjustment using `proc multtest` in SAS


```

proc rlogist data = hwl r = exchangeable semethod = zeger;
  nest strata psu unqid / psulev = 2;
  weight xwt;
  class sex agegrp urban income educ ethnic daily cpd ryob
        cohort country wave / nofreq dir = ascending;
  model wlnotice = sex agegrp urban income educ ethnic daily
        cpd ryob cohort country wave country*wave;
  /* Predicted marginals (adjusted probabilities) */
  predmarg country*wave;
  /* Tests difference in predicted marginals */
  pred_eff country = (1 0) * wave = (-1 1 0) / name = "Malaysia, W2 vs W1";
  pred_eff country = (1 0) * wave = (-1 0 1) / name = "Malaysia, W3 vs W1";
  pred_eff country = (1 0) * wave = (0 -1 1) / name = "Malaysia, W3 vs W2";
  pred_eff country = (0 1) * wave = (-1 1 0) / name = "Thailand, W2 vs W1";
  pred_eff country = (0 1) * wave = (-1 0 1) / name = "Thailand, W3 vs W1";
  pred_eff country = (0 1) * wave = (0 -1 1) / name = "Thailand, W3 vs W2";
  /* Relative change in support over time between countries */
  pred_eff country = (-1 1) * wave = (-1 1 0) / name = "TH W2-W1 vs MY W2-W1";
  pred_eff country = (-1 1) * wave = (-1 0 1) / name = "TH W3-W1 vs MY W3-W1";
  pred_eff country = (-1 1) * wave = (0 -1 1) / name = "TH W3-W2 vs MY W3-W2";
  /* Output estimates with p-values to dataset for proc multtest */
  output / PRMGCONS = default filename = contrasts replace;
run;

proc multtest
  pdata = contrasts(rename = (p_pmcon = RAW_P))
  out = padjust bon fdr;
run;

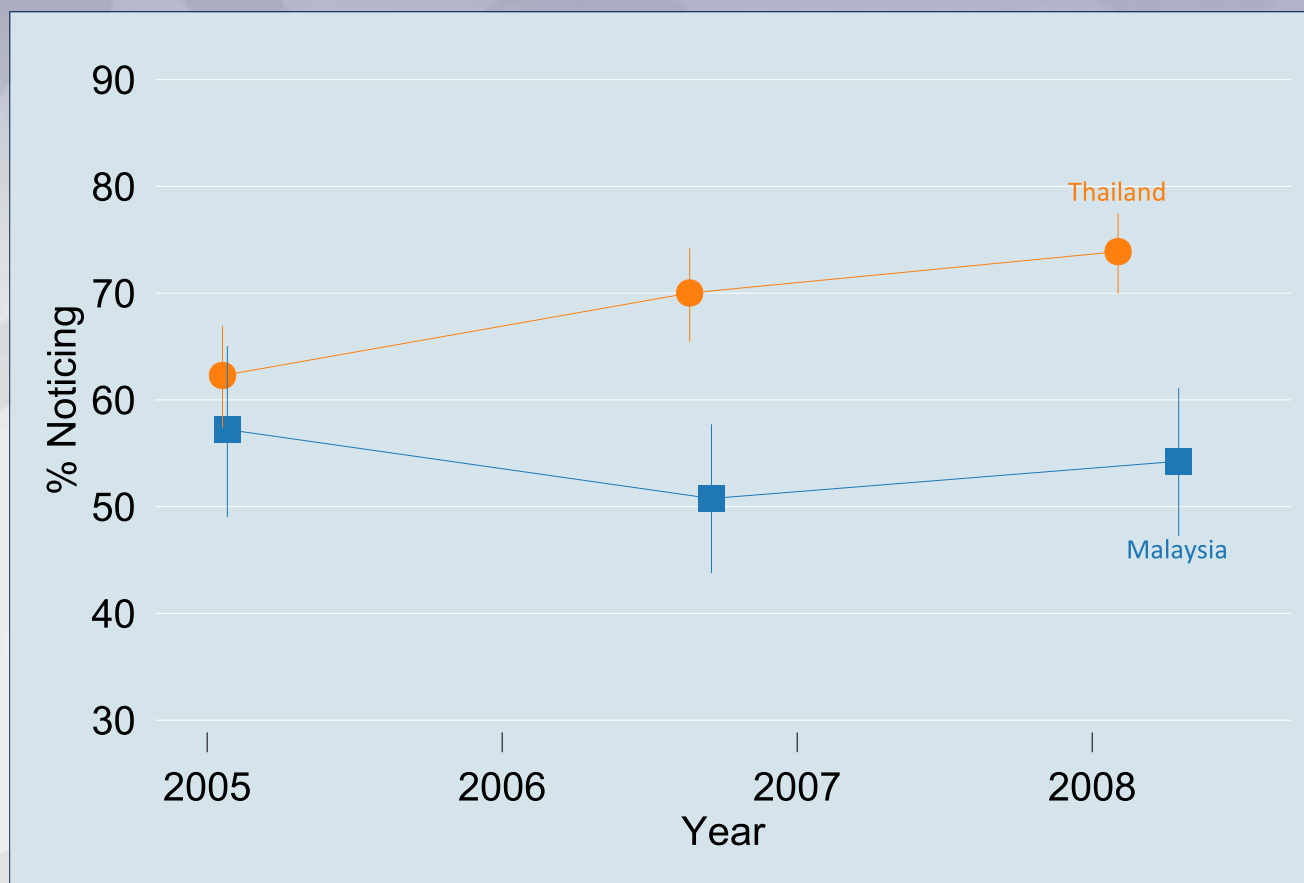
```

PRDEFFNO	PRMGCON	SEPMCON	T_PMCON	RAW_P	bon_p	fdr_p
Malaysia, W2 vs W1	-0.0507	0.0467	-1.0861	0.2828	1.00000	0.36354
Malaysia, W3 vs W1	0.0050	0.0486	0.1020	0.9192	1.00000	0.96281
Malaysia, W3 vs W2	0.0556	0.0410	1.3576	0.1808	1.00000	0.27121
Thailand, W2 vs W1	0.1264	0.0243	5.2105	0.0000	0.00003	0.00002
Thailand, W3 vs W1	0.1799	0.0267	6.7494	0.0000	0.00000	0.00000
Thailand, W3 vs W2	0.0536	0.0191	2.8012	0.0073	0.06542	0.01308
TH W2-W1 vs MY W2-W1	0.1771	0.0526	3.3635	0.0015	0.01351	0.00450
TH W3-W1 vs MY W3-W1	0.1750	0.0549	3.1885	0.0025	0.02243	0.00561
TH W3-W2 vs MY W3-W2	-0.0021	0.0443	-0.0469	0.9628	1.00000	0.96281

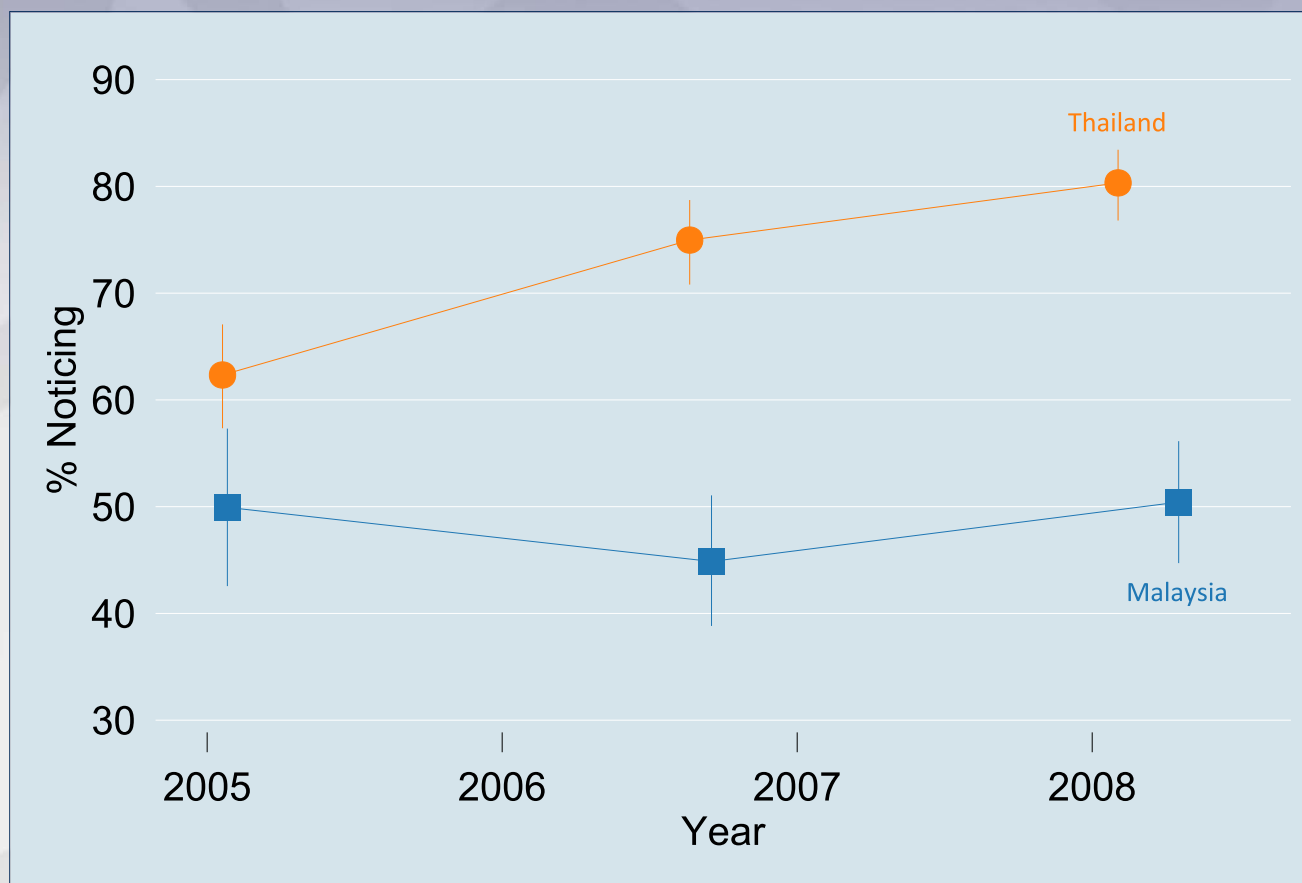
Connecting the Dots...

- Did warning label salience (noticing the health warning labels “often/very often”) increase significantly among smokers in Thailand following the introduction of new pictorial health warning labels in 2006?
- Was there any change in salience among Malaysian smokers where warning labels remained the same?
- Were any changes in salience among Thai smokers sustained over time?
- Used data from the first three waves of the ITC Southeast Asia (3,067 smokers from Thailand, 3,220 smokers from Malaysia, present in at least one wave)
- Binary GEE regression models used to
 - test hypotheses (controlling for other factors)
 - estimate the (adjusted) percentage of smokers noticing warning labels
 - account for the complex sampling design in the analysis **and** repeated measures (smokers could be present in 2 or all 3 waves)

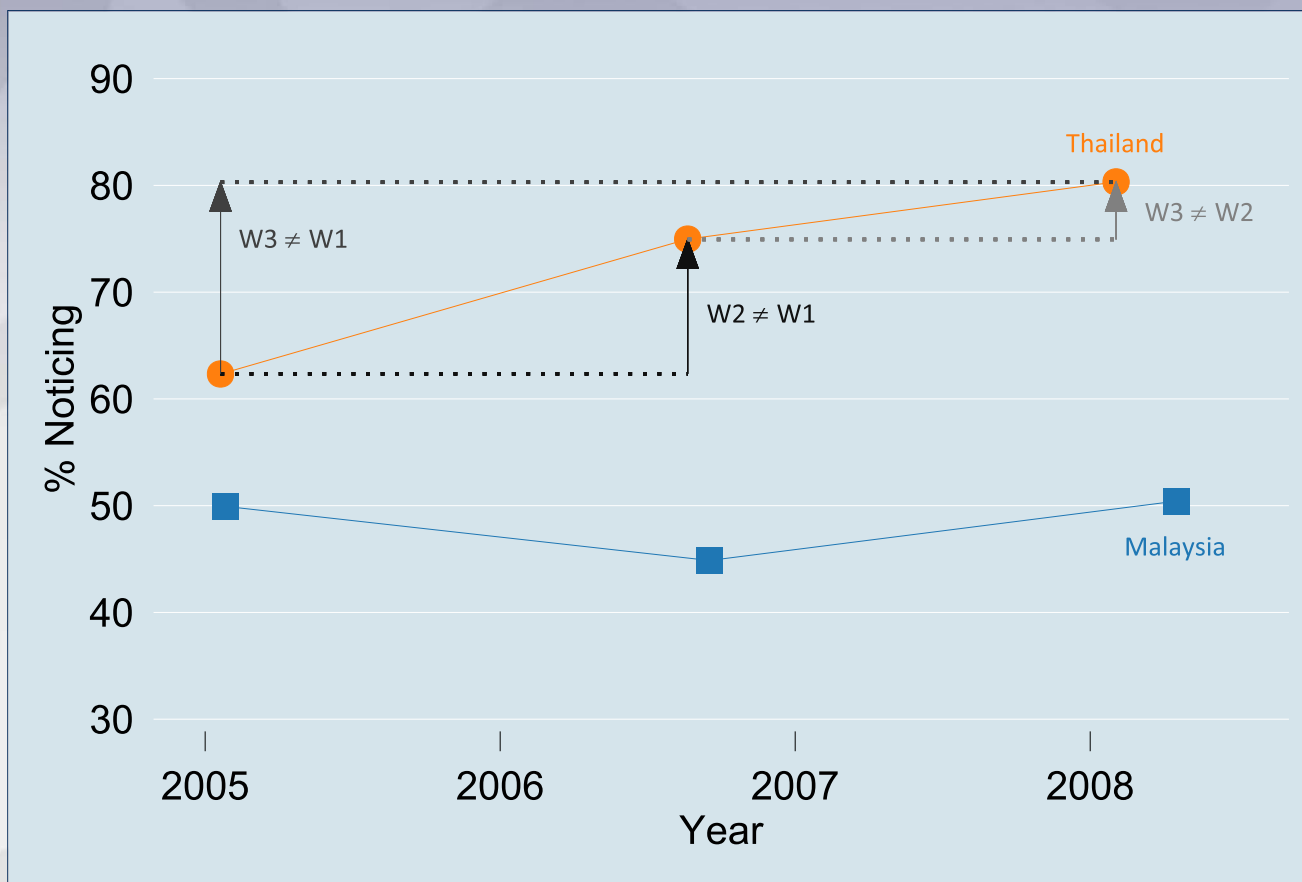
Unadjusted Estimates



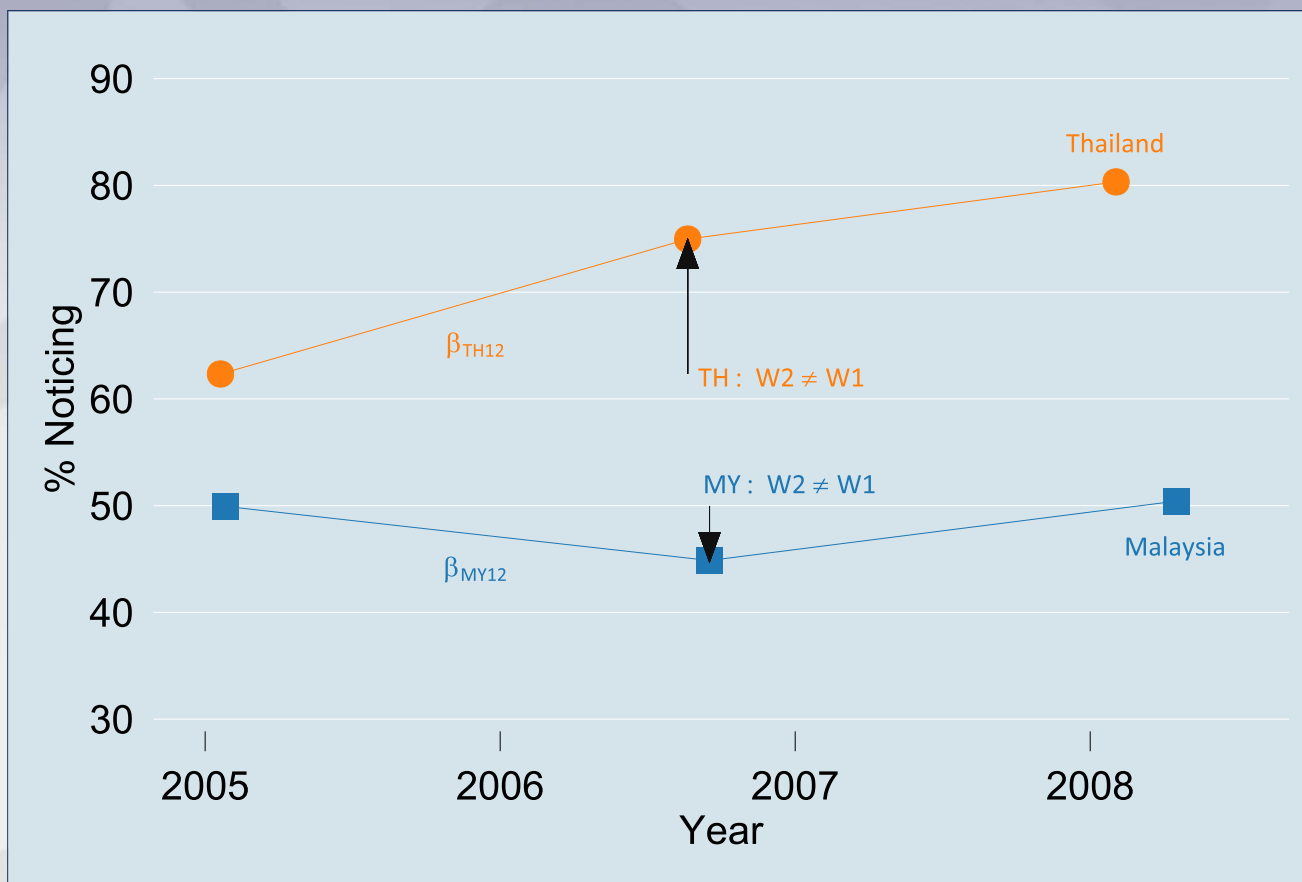
Adjusted Estimates



Hypothesis Tests



Other Possibilities



Statistical Tests

Contrast		% Diff.		p	OR	p
Thailand	1	62.3	—	—	—	—
	2	75.0	12.6	< 0.001 [†]	1.96	< 0.001 [†]
	3	80.3	5.4	0.007 [†]	1.41	0.007 [†]
Malaysia	1	49.9	—	—	—	—
	2	44.6	-5.1	0.283	0.80	0.282
	3	50.4	5.6	0.181	1.28	0.183
TH vs MY	W2 vs W1	—	17.7	0.002 [†]	2.46	0.001 [†]
	W3 vs W2	—	-0.2	0.963	1.10	0.658

[†] remains significant after controlling for multiple testing (false discovery rate).

Interpretation of Results

- Awareness of the health warning labels increased significantly among Thai smokers following introduction of the larger pictorial health warnings
- No changes in awareness were observed among Malaysian smokers
- The effect was sustained in Thailand by Wave 3 (about 3 years after the new health warning labels were introduced)
- Thailand's new pictorial health warnings have greater impact than the text-only warning labels they replaced and when refreshed, they help to reduce wear-out

Discussion

- Analyzing longitudinal data arising from complex survey designs such as those used by the ITC project is tricky!
- Typically, we wish to draw inferences about the population of smokers within ITC countries
- Population-averaged models, such as GEE models, are appropriate for analyzing longitudinal data in this situation
- However, it is still necessary to account for the complex sampling design in this situation to produce correct variance estimates and test statistics
- Although routines available in SAS for analyzing longitudinal data are suitable for stratified survey designs, for multi-stage designs where primary sampling units are clusters (e.g., villages or neighbourhoods), it is necessary to use SUDAAN to estimate GEE models

References & Suggested Reading

1. Hitchman SC, Driezen P, Logel C, Hammond D, Fong GT. (2014). Changes in effectiveness of cigarette health warnings over time in Canada and the United States, 2002–2011. *Nicotine & Tobacco Research*, 16: 536–543. doi: [10.1093/ntr/ntt196](https://doi.org/10.1093/ntr/ntt196).
2. Yong HH, Fong GT, Driezen P, Borland R, Quah ACK, Sirirassamee B, Hamann S, Omar M. (2013). Adult smokers' reactions to pictorial health warning labels on cigarette packs in Thailand and moderating effects of type of cigarette smoked: Findings from the International Tobacco Control Southeast Asia Survey. *Nicotine & Tobacco Research*, 15: 1339–1347. doi: [10.1093/ntr/nts241](https://doi.org/10.1093/ntr/nts241).
3. Ballinger GA. (2004). Using Generalized Estimating Equations for longitudinal data analysis. *Organizational Research Methods*, 7: 127–150. doi: [10.1177/1094428104263672](https://doi.org/10.1177/1094428104263672).
4. Gardiner JC, Luo Z, Roman LE. (2009). Fixed effects, random effects and GEE: What are the differences? *Statistics in Medicine*, 28: 221–239. doi: [10.1002/sim.3478](https://doi.org/10.1002/sim.3478).
5. Hanley JA, Negassa A, Edwardes MD, Forrester JE. (2003). Statistical analysis of correlated data using Generalized Estimating Equations: An orientation. *American Journal of Epidemiology*, 157: 364–375. doi: [10.1093/aje/kwf215](https://doi.org/10.1093/aje/kwf215).
6. Horton NJ, Lipsitz SR. (1999). Review of software to fit Generalized Estimating Equation regression models. *The American Statistician*, 53: 160–169. doi: [10.2307/2685737](https://doi.org/10.2307/2685737).
7. Zorn CJW. (2001). Generalized estimating equation models for correlated data: A review with applications. *American Journal of Political Science*, 45: 470–490. doi: [10.2307/2669353](https://doi.org/10.2307/2669353).

ITC Project Research Organizations



ITC Project Research Support



Core support provided by the Canadian Institutes of Health Research (MOP-115016) and the U.S. National Cancer Institute (P01 CA138389)