Exercise 6C: Genalized Linear Models and Extensions 2018 Spring

1 Background Knowledge

1.1 Quasi-Likelihood

- The maximum likelihood approach is based on the assumption of a specific model or data generation process (like Poisson process).
- When we have insufficient information about the data for use to specify a model for the data.
- but we can specify some features:
 - it is continuous or discrete,
 - how the mean/median is affected by extern variables,
 - how the variability of the response changes with the average,
 - whether the observations are independent,
 - whether the response distribution is skewed
- With these features we can develop analyses based on approximations to the likelihood.

Suppose

- We have a vector of responses, \mathbf{Y} , which are independent with mean μ and covariance matrix $\sigma^2 V(\mu)$
- μ is a function of covariates, x, and some regression parameters β . We can write this into $\mu(\beta)$.
- Typically, σ^2 is unknown, and has to be estimated.
- $V(\mu)$ is made up of known functions:

$$V(\mu) = \operatorname{diag}(V_1(\mu), \dots, V_n(\mu))$$

- We also assume that $V_i(\mu)$ only depends on μ_i .
- For a single component Y of Y:

$$U = u(\mu|Y) = \frac{Y - \mu}{\sigma^2 V(\mu)} \tag{1}$$

has several properties similar to log-likelihood derivative (i.e., the score),

$$\begin{array}{lcl} \mathbf{E}(U) & = & 0 \\ \mathbf{Var}(U) & = & 1/\left(\sigma^2V(\mu)\right) \\ -\mathbf{E}\left(\frac{\partial U}{\partial \mu}\right) & = & 1/\left(\sigma^2V(\mu)\right) \end{array}$$

• And the statistic:

$$Q(\mu|y) = \int_{y}^{\mu} u(y|y)dt = \int_{y}^{\mu} \frac{y-t}{\sigma^{2}V(t)}dt$$
 (2)

behaves like a log-likelihood function. We refer to this as log quasi-likelihood.

• The quasi-likelihood for the complete data:

$$Q(\mu|\mathbf{y}) = \sum Q(\mu_i|y_i)$$

• The quasi-deviance function for a single observation is

$$D(y|\mu) = -2\sigma^2 Q(\mu|y) = 2\int_y^\mu \frac{y-t}{V(t)} dt$$
 (3)

• The total deviance

$$D(\mathbf{y}|\mu) = \sum D(y_i|\mu_i)$$

only depends on μ and \mathbf{y} , but not σ^2 .

• The complete quasi-likelihood only depends multiplicately on σ^2 , so that it does not affect the MLEs of $\mu(\beta)$ and hence β .

Example 1 (Quasi-Gaussian Likelihood) When the variance function, $V(\mu) = 1$, then

$$U = \frac{Y - \mu}{\sigma^2}$$

so that the quasi-likelihood becomes

$$Q(\mu|y) = \int_{\mu}^{\mu} \frac{y-t}{\sigma^2} dt = -\frac{(Y-\mu)^2}{2}$$

which is the same as the likelihood for a normal distribution.

Example 2 (Quasi-Poisson Likelihood) When the variance function $V(\mu) = \mu$, then the quasi-score function:

$$U = \frac{Y - \mu}{\mu \sigma^2}$$

so the quasi-likelihood is

$$Q(\mu|y) = \int_{y}^{\mu} \frac{y - t}{t\sigma^{2}} = y \log \mu - \mu$$

which is the same as the likelihood for a Poisson distribution.

1.1.1 Other Quasi-likelihoods

1.2 Quasi-likelihood estimation

The quasi-score function is $\partial Q(\mu|\mathbf{y})/\partial \beta$, which is

$$\mathbf{U}(\beta) = \mathbf{D}^{\mathbf{T}} \mathbf{V}^{-1} (\mathbf{Y} - \mu) / \sigma^{2} = 0.$$
(4)

where

Table 1: Quasi-Likelihoods								
$V(\mu)$	$Q(\mu y)$	Distribution	Canonical parameter	Range restritions				
1	$-(y-\mu)^2/2$	Normal	μ	-				
μ	$y \log \mu - \mu$	Poisson	$\log \mu$	$\mu > 0$				
μ^2	$-y/\mu - \log \mu$	Gamma	$-1/\mu$	$\mu > 0, y \ge 0$				
μ^3	$-y/2\mu^2 + 1/\mu$	Inverse Gaussian	$-1/2\mu^{2}$	$\mu > 0, y \ge 0$				
μ^{ξ}	$\mu_{-\xi} \left(\frac{\mu y}{1-\xi} - \frac{\mu^2}{2-\xi} \right)$	-	$rac{1}{(1-\xi)\mu^{\xi-1}}$	$\mu>0,\xi\neq0,1,2$				
$\mu(1-\mu)$	$y\log(\mu/(1-\mu)) + \log(1-\mu)$	Binomial	$\log(\mu/(1-\mu))$	$0<\mu<1, 0\leq y\leq 1$				
$\mu + \mu^2/k$	$y \log\left(\frac{\mu}{k+\mu}\right) + k \log\left(\frac{k}{k+\mu}\right)$	Negative binomial	$\log\left(\frac{k}{k+\mu}\right)$	$\mu > 0, y \ge 0$				

- $\mathbf{D} \in \mathbb{R}^{n \times p}$ with $d_{ir} = \partial u_i / \partial \beta_r$
- $Cov(\mathbf{U}(\beta)) = -E(\partial \mathbf{U}(\beta))/\partial \beta$, and is

$$\mathbf{i}_{\beta} = \mathbf{D}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{D} / \sigma^{2}. \tag{5}$$

which is similar to the Fisher information for MLE.

• The asymptotic covariacne matrix of $\hat{\beta}$ is

$$Cov(\hat{\beta}) \approx \mathbf{i}_{\beta}^{-1} = \sigma^2 (\mathbf{D}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{D})^{-1}$$
(6)

• Estimation of σ^2 :

$$\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i} \frac{(Y_i - \mu_i)^2}{V_i(\hat{\mu}_i)} = X^2 / (n-p)$$
 (7)

1.3 Longitudinal Data Analysis

For longitudinal data, we can also use linear models with assumptions of correlations, when the number of observation per person, n_i , is small relative to the number of individuals m. However, for nonlinear discrete longitundinal data, we need to consider the different approach. There are three extensions of GLMs for longitudinal data:

- Marginal models
- Random Effects models
- Transition models

Marginal models

- Cross-sectional study
- Assumptions:
 - (1) The marginal expectation of the response, $E(Y_{ij}) = \mu_{ij}$, depends on the predictors, x_{ij} , by $h(\mu_{ij}) = x'_{ij}\beta$, where $h(\cdot)$ is a known link function.
 - (2) The marginal variance $Var(Y_{ij}) = v(\mu_{ij})$ where $v(\cdot)$ is a known variance function and ϕ is a scale parameter.

- (3) The correlation $Corr(Y_{ij}, Y_{ik}) = \rho(\mu_{ij}, \mu_{ik}; \alpha)$ where $\rho(\cdot)$ is a known function.
- Marginal models are natural analogues for correlated data of GLMs for independent data.

1.4 Generalized Estimating Equations (GEEs)

Generalized estimating equations (GEEs) are a specific family of marginal models, which are used to model correlated data from

- Longitudinal/repeated measures studies: for same subjects, same measures, successive times. And the succesive measurements are expected to be correlated.
- Clustered/multi-level studies

1.4.1 Notations

For a set of repeated measurements y_{ij} , where i = 1, ..., N for each subject; while $j = 1, ..., n_i$ be the times for subject i.

Similarly, for clustered data y_{ij} , where i = 1, ..., N denotes clusters; while $j = 1, ..., n_i$ denotes measurements within cluster i.

In a normal linear model, for unit i

$$E(y_i) = \mu_i = X_i \beta$$

$$y_i \sim N(\mu_i, V_i)$$

where

- $X_i \in \mathbb{R}^{n \times p}$ is the design matrix;
- $\beta \in \mathbb{R}^p$ is the parameter vector;
- $V_i \in \mathbb{R}^{n_i \times n_i}$ is the variance-covariance matrix (e.g., $V_i = \sigma_i^2 I$ if measurements are independent).

For all units:

$$E(\mathbf{y}) = \mu = \mathbf{X}\beta$$
$$\mathbf{y} \sim N(\mu, \mathbf{V})$$

where

$$\mu = (\mu_1, \dots, \mu_N)^T, \mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_N)^T, \mathbf{V} = \operatorname{diag}(\mathbf{V}_1, \dots, \mathbf{V}_N).$$

Then we need to estimate β and \mathbf{V} in the model.

Similarly, we use the log-likelihood function:

$$l = (\mathbf{v} - \mu)^{\mathbf{T}} \mathbf{V}^{-1} (\mathbf{v} - \mu)$$

and the score function:

$$U(\beta) = \frac{\partial}{\partial \beta} \mathbf{X^T} \mathbf{V^{-1}} (\mathbf{y} - \mu)$$

which can be solved using a set of score equations:

$$\mathbf{X_i^T}\mathbf{V^{-1}}(\mathbf{y_i} - \mathbf{X_i}\boldsymbol{\beta}) = \mathbf{0}$$

1.4.2 Marginal models for generlized linear model

Similarly,

$$E(Y_{ij}) = \mu_{ij}, g(\mu_{ij}) = \eta_{ij} = x_i \beta$$

and the score function becomes

$$U(\beta) = \mathbf{D_i^T V_i^{-1}}(\mathbf{y_i} - \mu_i) = \mathbf{0}$$

where \mathbf{D}_i is a matrix of derivatives with elements:

$$\frac{\partial \mu_i}{\partial \beta_k} = \frac{\partial \mu_i}{\partial \eta_k} x_{ik}$$

and V_i is diagonal with elements $Var(Y_{ij})$.

1.4.3 GEEs

Since Y_{ij} 's are not necessarily independent, if $\mathbf{R}_i \in \mathbb{R}^{n_i \times n_i}$ is the correlation matrix for cluster i, then the variance-covariance matrix \mathbf{V} can be rewritten as:

$$\mathbf{V}_i = \mathbf{A_i^{1/2}} \mathbf{R_i} \mathbf{A_i^{1/2}}$$

where \mathbf{A}_i is the diagonal matrix with elements $\text{Var}(Y_{ij})$. And thus

$$U(\beta) = \mathbf{D_i^T V_i^{-1}}(\mathbf{y_i} - \mu_i) = 0$$

where $V_i = A_i^{1/2} R_i A_i^{1/2}$.

- \mathbf{D}_i is the matrix of derivatives: $\frac{\partial \mu_i}{\partial \beta_k}$;
- V_i is the "working" variance-covariance matrix of Y_i ;
- $\mathbf{A}_i = \operatorname{diag}\{\operatorname{var}(Y_{ik})\}$
- \mathbf{R}_i is the correlation matrix for Y_i
- ϕ is an overdispersion parameter.

1.5 Random Effects Models

The previous model, either linear or generalized linear, can be used to estimate the "fixed" effects, which consist of specific and repeatable categories/variables that are representative of an entire population (e.g., species, gender, age). In a longitudinal study with repeated mesures, and in studies using hierarchical (nested) sampling, it is also possible to estimate effects associated with individuals sampled at random from the population of interest. These are "random" effects which convey information about the degree that individuals in a population differ but not how or why they differ.

Then how to differentiate fixed and random effects?

- Fixed effects can capture informations that are beyond the current analysis (a species of tree)
- Random effects contain only the information that are not beyond the current analysis (a group of tree within an observation)

- Fixed effects influence the *mean* of the response.
- Random effects only infludence the *variance* of the response.

Then why mixed-effects models? Mixed-effects models are particularly useful to deal with potential pseudoreplication and unbalanced designs. Including random effects can also account for variation that can mask patterns if we only consider fixed effects.

1.5.1 General specification of random-effects GLM

(1) Given U_i , the response Y_{i1}, \ldots, Y_{in_i} are mutually independent and

$$f(y_{ij}|U_i) = \exp\left[\frac{y_{ij}\theta_{ij} - b(\theta_{ij})}{\phi} + c(y_{ij}, \phi)\right]$$

(2) The conditional moments:

$$\mu_{ij} = E(Y_{ij}|U_i) = b'(\theta_{ij})$$

$$v_{ij} = Var(Y_{ij}|U_i) = b''(\theta_{ij})\phi$$

which satisfy:

$$h(\mu_{ij}) = x'_{ij}\beta^* + d'_{ij}U_i$$

$$v_{ij} = v(\mu_{ij})\phi$$

where $h(\cdot)$ and $v(\cdot)$ are known link and variance functions, and d_{ij} is a subset of x_{ij} .

(3) The random effects, U_i , i = 1, ..., m are mutually independent with a common underlying multivarite distribution, F.

1.5.2 Basic underlying random effects model

- There is natural heterogeneity across individuals in their regression coefficients (which can be represented by a probability distribution)
- Correlation among observations for one person arises from their sharing of unobservable variable, U_i .
- Also known as *latent variable* model

Contrast to the marginal models, the random effects model is especially useful when the objective is to make inference about individuals rather than the population average.

Example 3 (Correlation from random effects)

$$Y_{ij} = \mu + u_i + \epsilon_{ij}$$

where

- Y_{ij} : response for unit i in repeated j;
- μ : average value for population.
- u_i : random effect with $u_i \sim N(0, \sigma_u^2)$

• ϵ_{ij} : error with $\epsilon_{ij} \sim N(0, \sigma_e^2)$

therefore,

- $E(Y_{ij}) = \mu$
- $Var(Y_{ij}) = \sigma_u^2 + \sigma_e^2$
- The covariance matrix:

$$cov(Y_{ij}, Y_{km}) = \begin{cases} \sigma_u^2 & \text{if } i = k \\ 0 & \text{otherwise} \end{cases}$$

So V_i is exchangeable with elements $\frac{\sigma_u^2}{\sigma_u^2 + \sigma_e^2}$, which is called intra-class correlation coefficient (ICC).

2 Exercises

- 1. (20 points) Write down the quasi-score equation for
 - (1) Quasi-Poisson distribution
 - (2) Quasi-Binomial distribution
 - (3) Quasi-Gamma distribution
 - (4) Quasi-Inverse Gaussian distribution
 - (5) Quasi-Negative binomial distribution.
- 2. (20 points) Assume that the random variable Y is Poisson distributed with probability mass function

$$P(Y = y|\lambda) = \frac{\lambda^y}{y!} \exp(-\lambda)$$

(1) Show that the distribution of Y belongs to the exponential distribution family. That is show that the function can be rewritten as the form

$$\exp\{[\theta y - b(\theta)]/a(\phi) + c(y,\phi)\},\$$

and determine $\theta, b(\theta), a(\phi)$ and $c(y, \phi)$.

Assume that Y_1, \ldots, Y_n are independent with the Poisson distribution, and let $\mu_i = E(Y_i), i = 1, \ldots, n$

- (2) Explain what we mean by a generalized linear model (GLM) for Y_i with link function g, and determine the canonical link function.
- (3) Derive an expression for the log-likelihood function $L(\mu; \mathbf{y})$ where $\mathbf{y} = (y_1, \dots, y_n)^T$ is the observed value of $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ and $\mu = (\mu_1, \dots, \mu_n)^T$.
- (4) Explain what we mean by a saturated model and determine the maximum of $L(\mu; \mathbf{y})$ for the saturated model.
- (5) Explain what we mean by the deviance $D(\mathbf{y}; \hat{\mu})$ of a Poisson GLM, find an expression for the deviance, and discuss how it may be used.

Table 2: Ten years of deaths from coronary diseases								
Age	Smokers			Non-smokers				
group	Deaths Person-years		-	Deaths	Person-years			
35-44	32	52407		2	18790			
45-54	104 43248			12	10673			
55-64	206	28612		28	5710			
65-74	186	12663		28	2585			
75-84	102	5317		31	1462			

3. (10 points) Mittlbock and Heinzl (2001) compare Poisson and logistic regression models for data in which the event rate is small so that the Poisson distribution provides a reasonable approximation to the Binomial distribution. An example is the number of deaths from coronary heart disease among British doctors.

We can fit the model $Y_i \sim \text{Poisson}(\text{deaths}_i)$ with the following equation:

$$\log(\text{deaths}_i) = \log(\text{personyears}_i) + \beta_1 + \beta_2 \text{smoke}_i + \beta_3 \text{agecat}_i + \beta_4 \text{agesq}_i + \beta_5 \text{smkage}_i$$

An alternative is $Y_i \sim \text{Bin}(\text{personyears}_i, \pi_i)$ with

$$logit(\pi_i) = \beta_1 + \beta_2 smoke_i + \beta_3 agecat_i + \beta_4 agesq_i + \beta_5 smkage_i$$
.

Another version is based on a Bernoulli distribution $Z_i \sim Bernoulli(\pi_i)$ for each doctor in group i with

$$Z_{j} = \begin{cases} 1 & j = 1, \dots, \text{deaths}_{i} \\ 0 & j = \text{deaths}_{i} + 1, \dots, \text{personyears}_{i} \end{cases}$$

and

$$logit(\pi_i) = \beta_1 + \beta_2 smoke_i + \beta_3 agecat_i + \beta_4 agesq_i + \beta_5 smkage_i$$
.

- (1) (5 points) Fit all three models. Verify that the β estimates are very similar.
- (2) (5 points) Calculate the statistics D, X^2 and pseudo R^2 for all three models. Notice that the pseudo R^2 is much smaller for the Bernoulli model. This is probably due to the reason that the Poisson and Binomial models are estimating the probability of deaths for each group (which is relatively easy) whereas the Bernoulli model is estimating the probability of deaths for an individual (which is much more difficult).
- 4. (30 points) Twenty-four patients with stroke are randomized to three different treatments:
 - A: new OT intervention;
 - B: special stroke unit in the same hospital;
 - C: usual care in different hospital.

Each group has 8 patients.

The primary outcome is the meaurement of functional ability - Barthel index, which was measured weekly for 8 weeks.

```
stroke <- read.table("stroke.dat", header=TRUE)</pre>
stroke.long <- reshape(stroke, idvar=c("Subject", "Group"),</pre>
```

varying=3:10, timevar="Week", direction="long")
rownames(stroke.long) <- NULL
names(stroke.long)[4] <- "Ability"
stroke.long\$Group <- factor(stroke.long\$Group)
stroke.long\$Subject <- factor(stroke.long\$Subject)</pre>

(1) (5 points) Pooled analysis ignoring correlation within patients

$$Y_{ijk} = \alpha_j + \beta_j k + e_{ijk}$$

where j denotes for groups, k for times, and i for subjects. Here we use different intercepts and different slopes for groups. Assume that all Y_{ijk} are independent and of the same variance (i.e. ignoring the correlation between observations). Use multiple regression to compare α_j 's and β_j 's. Here we use the interaction term to model the different slopes between different groups.

(2) (5 points) Data reduction Fit a linear model for each patient:

$$Y_{ijk} = \alpha_{ij} + \beta_{ij}k + e_{ijk}$$

This also assume independence and constant variance.

- (A) Use simple linear regression to estimate α_{ij} and β_{ij} .
- (B) Perform ANOVA using estimates $\hat{\alpha}_{ij}$ as the observations and groups as levels of a factor in order to compare α_j 's.
- (C) Similarly compare β_j 's using $\hat{\beta}_{ij}$ as the observations.
- (3) (5 points) Repeated-measures analysis using various variance-covariance structures Fit

$$Y_{ijk} = \alpha_i + \beta_i k + e_{ijk}$$

where

• α_i and β_i are the parameters of interest.

Assume normality for e_{ijk} but try various forms for variance-covariance matrix.

For the stroke data, choose the **auto-regression structure** (e.g., AR(1)) as the appropriate model. And then use GEEs to fit the models.

Note: In R, you need to apply the geepack::geeglm() function.

(4) (5 points) Mixed/Random-effects model Use the model

$$Y_{ijk} = (\alpha_j + a_{ij}) + (\beta_j + b_{ij})k + e_{ijk}$$

where

- α_i and β_i are fixed effects for groups.
- $a_{ij} \sim N(0, \sigma_a^2), b_{ij} \sim N(0, \sigma_b^2)$ and $e_{ijk} \sim N(0, \sigma_e^2)$ are random effects and all are independent.
- Fit this mixed model and use estimates of fixed effects to compare α_i 's and β_i 's.

In R, you need to apply the nlme::lme() function.

(5) (10 points) Make a conclusion based on the above models and compare the advantages and disadvantage of the above models.

Table 3: Data of stroke recovery

		Week							
Subject	Group	1	2	3	4	5	6	7	8
1	A	45	45	45	45	80	80	80	90
2	A	20	25	25	25	30	35	30	50
3	A	50	50	55	70	70	75	90	90
4	A	25	25	35	40	60	60	70	80
5	A	100	100	100	100	100	100	100	100
6	A	20	20	30	50	50	60	85	95
7	A	30	35	35	40	50	60	75	85
8	A	30	35	45	50	55	65	65	70
9	В	40	55	60	70	80	85	90	90
10	В	65	65	70	70	80	80	80	80
11	В	30	30	40	45	65	85	85	85
12	В	25	35	35	35	40	45	45	45
13	В	45	45	80	80	80	80	80	80
14	В	15	15	10	10	10	20	20	20
15	В	35	35	35	45	45	45	50	50
16	В	40	40	40	55	55	55	60	65
17	С	20	20	30	30	30	30	30	30
18	С	35	35	35	40	40	40	40	40
19	С	35	35	35	40	40	40	45	45
20	С	45	65	65	65	80	85	95	100
21	С	45	65	70	90	90	95	95	100
22	С	25	30	30	35	40	40	40	40
23	С	25	25	30	30	30	30	35	40
24	С	15	35	35	35	40	50	65	65

5. (20 points) Assume that $U_i \sim N(0, \sigma^2)$. Given $U_i = u_i$, the binary random variables Y_{i1}, \ldots, Y_{in_i} are independent with

$$P(Y_{ij} = 1 | U_i = u_i) = 1 - P(Y_{ij} = 0 | U_i = u_i) = \Phi(\beta_0 + \beta_1 x_{ij} + u_i)$$

where $\Phi(\cdot)$ is the CDF for standard normal distribution, and x_{ij} 's are known.

(1) What is this model called? Describe one or more situations where such a model can be useful.

A marginal model for Y_{ij} 's is given by:

$$P(Y_{ij} = 1) = 1 - P(Y_{ij} = 0) = \Phi(\gamma_0 + \gamma_1 x_{ij})$$

- (2) Show how the parameters γ_0 and γ_1 may be expressed in terms of β_0, β_1 and σ^2 .
- (3) Comment on the marginal model and the random-effects models for such a clustered binary data, based on the above questions.
- 6. (20 points) As we know, the saturated model has a separated parameter μ_i with no contraint. Let $\hat{\theta}_i$ and $\tilde{\theta}_i$ denote the parameter under model ω and Ω (saturated model), respectively. The likelihood ratio test (LRT) criterion to compare the two models in the exponential family has

the form:

$$-2\log\lambda = -2\log\frac{L(\hat{\theta})}{L(\tilde{\theta})} = 2\left[l_{\Omega}(\tilde{\theta})\right] = 2\sum_{i=1}^{n}\frac{y_{i}(\tilde{\theta}_{i} - \hat{\theta}_{i}) - (b(\tilde{\theta}_{i}) - b(\hat{\theta}_{i}))}{a_{i}(\phi)}$$

assuming that $a_i(\phi) = \phi/w_i(w_i = 1)$ for known prior weight w_i .

In GLM, the deviance for model ω can be defined as

$$D(\omega) = 2a(\phi) \left[l_{\Omega}(\tilde{\theta}) - l_{\omega}(\hat{\theta}) \right] = 2 \sum_{i=1}^{n} \left[y_i(\tilde{\theta}_i - \hat{\theta}_i) - (b(\tilde{\theta}) - b(\hat{\theta})) \right].$$

Thus

$$-2\log\lambda = D(\omega)/a(\phi)$$

is called the scaled deviance.

In order to compare two nested models:

$$-2\log\lambda = -2\log\left[\frac{L_{\omega_1(\theta_1)}}{L_{\omega_2}(\theta_2)}\right] = 2\{[l_{\Omega}(\theta) - l_{\omega_1}(\theta_1)] - [l_{\Omega}(\theta) - l_{\omega_2}(\theta_2)]\} = \frac{D(\omega_1) - D(\omega_2)}{a(\phi)} \overset{n\to\infty}{\sim} \chi^2_{p_2-p_1}$$

Here the scale parameter $a(\phi)$ is either known or estimated using the larger model, ω_2 .

Write down the deviance for

- (1) Normal distribution.
- (2) Poisson distribution.
- (3) Binomial distribution.
- 7. (20 points) The measurement of *left ventricular volume* of the heart is important for studies of cardiac physiology and clinical management of patients with heart disease. An indirect way of measuring the volume, y, involves a measurement called *parallel conductance volume*, x. Boltwood et al. (1989) found an approximately linear association between y and x in a study of dogs under various "load" conditions. The results, reported by Glantz and Slinker (1990), are shown in the following table.

Table 4: Measurements of left ventricular volume and parallel conductance volume

		Condition									
Dog		1	2	3	4	5	6	7	8		
1	у	81.70	84.30	72.80	71.70	76.70	75.80	77.30	86.30		
	X	54.30	62.00	62.30	47.30	53.60	38.00	54.20	54.00		
2	y	105.00	113.60	108.70	83.90	89.00	86.10	88.70	117.60		
	X	81.50	80.80	74.50	71.90	79.50	73.00	74.70	88.60		
3	y	95.50	95.70	84.00	85.80	98.80	106.20	106.40	115.00		
	X	65.00	68.30	67.90	61.00	66.00	81.80	71.40	96.00		
4	y	113.10	116.50	100.80	101.50	120.80	95.00	91.90	94.00		
	X	87.50	93.60	70.40	66.10	101.40	57.00	82.50	80.90		
5	y	99.50	99.20	106.10	85.20	106.30	84.60	92.10	101.20		
	X	79.40	82.50	87.90	66.40	68.40	59.50	58.50	69.20		

- (1) (EDA) Conduct a exploratory analysis of these data.
- (2) (**Pooled analysis**) Let (Y_{jk}, x_{jk}) denote the k-th measurement (k = 1, ..., 8) on dog j (j = 1, ..., 5). Fit the linear model:

$$E(Y_{jk}) = \mu = \alpha + \beta x_{jk}, Y \sim N(\mu, \sigma^2).$$

- (3) (Data reduction analysis) Fit a linear model based on the data reduction approach.
- (4) (Random-effects model) Fit a suitable random effects model by specifying the fixed effects and random effects.
- (5) (**GEE**) Fit a clustered model using a GEE.
- (6) Compare the results obtained from each approach. Which method(s) do you prefer? Why?