# On the Joint Analysis of Longitudinal Responses and Early Discontinuation in Randomized Trials

François Vandenhende[1,2] and Philippe Lambert[2]

[1] Statistics and Information Sciences, Eli Lilly & Company,

B-1348 Mont-Saint-Guibert, Belgium.

[2] Institute of Statistics UCL, Belgium.

June 17, 2002

Summary. Our focus is on the joint analysis of longitudinal non-normal responses and early discontinuation in (pre)-clinical trials. Separate models are fitted to the two series (response and discontinuation) to account for covariate and time e ects. The serial dependence and the dependence between response and drop-out are also modelled. This is done using particular dependence functions, called copulas. Copulas are used to create a joint distribution with given marginal distributions. Applications are given for the analysis of heart rate/morbidity in toxicology and pain severity/intake of rescue medications in a trial on migraine. Using copulas, the level of dependence between two variables remains invariant to changes in the marginal distribution of either variable. This proves interesting in modelling the association in a longitudinal setting when responses change over time.

*email:* francois@lilly.com

1

Key words: Copula; Drop-out; Repeated Measures; Continuous or Ordinal Data; Likelihood.

## 1. Introduction

In many longitudinal studies, subjects drop out of the trial before completion. This produces unbalanced sets of data with unequal numbers of observations per subject. In addition to the longitudinal response, the drop-out rate can also be interesting to analyze. This was the case in a recent migraine trial in which drop-out occurred after the intake of a rescue medication (RM). Both the pain severity profile and the rate of RM intake can be analyzed to provide evidence of a potential treatment e ect. Standard statistical software is available to analyze unbalanced repeated and survival data separately. But when the drop-out process is not independent of the response - e.g. when patients tend to take a RM as their pain remains severe - a separate analysis of the two variables can bias the results.

Rubin (1976) introduced and Little and Rubin (1987) developed a taxonomy to distinguish several types of missing data according to their association to the response. The non-response process is termed *missing completely at random* (MCAR) when it is independent of responses. It is *missing at random* (MAR) when conditionally independent of the unobserved response(s) given the observed data. When neither MCAR nor MAR hold, it is termed *missing not at random* (MNAR). The analysis of repeated measures with non random drop-outs has been discussed by various authors. Likelihood methods were used by Diggle and Kenward (1994) for continuous responses and by Molenberghs, Kenward and Lesa re (1997) for discrete ordinal responses. In the two cases, the joint density $f^{Y,D}$ of the response vector Y and drop-out

time $D$ is factored as $f^Y f^{D/Y}$. These classes of models are known as *selection models*. The conditional distribution of ($D/$Y) is typically modelled as an event-time regression including Y as covariate. A distinction between MCAR, MAR and MNAR processes is made based upon the responses which have to be included in the modelling of $f^{D/Y}$.

Another factoring is $f^{Y/D} f^D$. This is known as *pattern mixture* models. This approach stratifies the data over the different patterns (e.g. times) of drop-out and then applies a different model to each pattern. An application is given by Little (1993). A taxonomy equivalent to that of Rubin (1976) was developed by Molenberghs, Michiels, Kenward and Diggle (1998) in the framework of pattern mixture models.

Both selection and pattern mixture models have merit in either the longitudinal or the survival settings but the interpretation of parameters is different in the two models. A marginal and a conditional model are explicitly available in the two cases. The conditional model applies to the drop-out process in a selection model and to the response in a pattern mixture specification. The transformation from one probability model to the other is generally not straightforward (Little, 1994).

When the two conditional models $f^{Y/D}$ and $f^{D/Y}$ are of direct interest, Molenberghs, Michiels and Kenward (1998) have developed a pseudo-likelihood method for their joint analysis. In this paper, we focus on the specification of marginal models $f^Y$ and $f^D$ for the response and drop-out indicator, respectively. We develop a likelihood method that permits a valid joint inference on these models when drop-outs are not MCAR. To do so, we create a joint distribution for Y and $D$ in which the dependence between the

marginal distributions is modelled using *copulas*.

A copula $C(\cdot)$ is a multivariate distribution on $[0,1]^J$, where $J$ is the variable dimension. Sklar (1959) showed that for any set of random variables $Z_1, ..., Z_J$ with univariate marginal distributions $F_1(z_j)$, the construction

$$F_J(z_1, ..., z_J) = C[F_1(z_1), ..., F_1(z_J)], \tag{1}$$

generates a $J$-variate joint distribution for the $Z_j$. Hence, we propose to use parametric families of copulas to generate a joint density $f^{Y,D}$ with marginal densities $f^Y$ and $f^D$. Given $f^{Y,D}$, a likelihood-based analysis is possible to jointly estimate the parameters of the margins and of the copula which specify the dependence structure. Since the dependence embodied in the copula does not depent on the margins, its parameters could also be estimated nonparametrically (Shih and Louis, 1995; Genest et al., 1995).

As a joint distribution is obtained whatever the marginal distributions $F_1(z_j)$ involved in (1), copula models can thus evenly be applied to continuous or to non-continuous responses. Oakes (1994), Shih and Louis (1995) and Wang and Wells (2000) use copulas to model multivariate survival data. Meester and MacKay (1994) gave an application to the modelling of clustered categorical data. Trégouèt, Ducimetière, Bocquet, Visvikis, Soubrier and Tiret (1999) presented a copula model for the analysis of familial binary data. In this paper, we consider longitudinal data of either continuous or discrete (ordinal) type and we also use copulas to model the dependence between repeated responses. The techniques proposed by Lambert and Vandenhende (2002) and by Vandenhende and Lambert (2000) are used to model the dependence in continuous and ordinal data respectively.

The layout of the paper is as follows. Notations are introduced in Section 2. The joint likelihood for repeated responses and drop-outs is given in Section 3. Copula models for repeated measures are defined in Section 4 and models for the dependence between response and drop-out are in Section 5. In Section 6, we present the copulas that will be used in the examples and review their dependence properties. Two examples are given in Section 7. A continuous response (heart rate) is modelled in a toxicology study and a discrete ordinal response (pain severity score) is considered in a clinical trial on migraine. A discussion of key results is then made in Section 8.

## 2. Some Notations

We consider longitudinal studies with the same number of measurements scheduled for all subjects. For simplicity, notations are introduced without subscript for the subject's indicator. Let the random vector Y be the complete set of measurements on a subject and let O be the associated *response indicator*. For a particular realization $(y, o)$, the elements of o take the value 1 when the corresponding element of y is observed and a value 0, when missing. This paper does not consider intermittent but only monotone patterns of missing data (i.e., drop-outs). In this situation, the information contained in the vector O can be replaced by a random variable $D$ indexing the first time $t_d$ at which no response is available for the subject. We shall set $d = \ +1$ for completers. A series of covariates is recorded together with the data. Parametric models indexed by a vector of parameters will be considered to relate covariates to the responses and to drop-out hazards. The vector $^T$ is partitioned as ($_Y^T$, $_D^T$, $^T$) where $_Y$ and $_D$ are for the marginal models of the response and drop-out respectively. The dependence

5

model between Y and $D$ is given by a copula $K(\cdot;\ )$ with parameter . The

vector $\bar{\gamma}$ is further partitioned into ( $^T$, $^T$) where is the set of parame-

ters for the marginal models at each time. The dependence model between

repeated responses is given by a copula $Q(\cdot;\ )$ indexed by . Distributions

and densities (pmf for discrete responses) are denoted as $F_j(\cdot;\ )$ and $f_j(\cdot;\ )$

respectively, where the subscript $j$ gives the dimension.

## 3.  Likelihood for a Selection Model

Let us consider the realizations (y, o) for any subject dropping at the $d^{th}$ time.

In a selection model, the contribution to the likelihood $f_d(y_1, ..., y_{d-1}, d;\ )$

for this subject can be written as (Diggle and Kenward, 1994)

$$f_{d-1}(y_1, ..., y_{d-1};\ \gamma)\ \prod_{j=1}^{d-1} Pr(O_j = 1/O_{j-1} = 1, y_1, ..., y_j;\ )$$

$$Pr(O_d = 0/O_{d-1} = 1, y_1, ..., y_d;\ )g_1(y_d/y_1, ..., y_{d-1};\ \gamma)dy_d. \tag{2}$$

The first factor in (2) is the joint density of all observed responses $y_1, ..., y_{d-1}$

for the subject; the next factor is the joint probability of being observed up

to time $t_{d-1}$ conditional on the historical responses; and the last factor ex-

pressed as an integral is the expected (conditional) drop-out hazard at time

$t_d$ under the full conditional density $g_1(y_d/y_1, ..., y_{d-1})$ for the non-observed

response $Y_d$.

For subjects who complete the study, $d = \ + 1$ and (2) reduces to

$$f\ (y_1, ..., y\ ;\ \gamma)\ \prod_{j=1} Pr(O_j = 1/O_{j-1} = 1, y_1, ..., y_j;\ ).$$

Under MCAR mechanisms, the drop-out hazard is totally independent of

y and (2) can be rewritten as

$$f_{d-1}(y_1, ..., y_{d-1}; \ _Y) \prod_{j=1}^{d} \Pr(O_j = o_j / O_{j-1} = 1; \ _D).$$ (3)

The density of the observed responses and the probability model for the drop-out hazard clearly separate in (3). Parameter vectors $_Y$ and $_D$ are distinct and a separate estimation of the two models (e.g. using maximum likelihood) produces unbiased estimates.

When the drop-out process is MAR, Equation (2) becomes

$$f_{d-1}(y_1, ..., y_{d-1}; \ _Y) \prod_{j=1}^{d} \Pr(O_j = o_j / O_{j-1} = 1, y_1, ..., y_{j-1}; \ _Y, \ _D, \ ).$$ (4)

Little and Rubin (1987) have suggested that, under MAR mechanisms, the drop-out process could be *ignorable* for a likelihood-based analysis of the response model. The ignorability condition requires that the set of parameters for the models $f^Y$ and $f^{D/Y}$ be distinct (the assumption of separability). In our case, we do not specify a parametric model for $f^{D/Y}$ directly. Instead, it is a consequence of the chosen marginal models $f^D$ and $f^Y$ and of the copula used to relate them. The complete vector of parameters ( $_Y^T$, $_D^T$, $^T$) is involved in the modelling of $f^{D/Y}$. The assumption of separability between $_Y$ and ( $_D^T$, $^T$) is not always verified in our approach so that ignorability does not necessary hold for all dependence models. As a consequence, maximum likelihood estimators of $_Y$ can be di erent between an identical response model specified in (3) and in (4).

## 4. Dependence between Repeated Responses

We start by specifying a common parametric family $F_1(y_j; \ )$ for the distribution of $Y_j$ at each time $t_j$ ($j = 1, ..., d$). For any subject dropping at $t_d$, the

joint distribution of the $d-1$ successive observations is then modelled as in (1) using a copula $Q(\cdot;\ )$ indexed by a vector of parameters . This yields

$$F_{d-1}(y_1,...,y_{d-1};\ \gamma) = Q[F_1(y_1;\ ),...,F_1(y_{d-1};\ );\ ]. \qquad (5)$$

When the response is absolutely continuous, the joint density is obtained by derivation of (5) with respect to $Y_1,...,Y_{d-1}$. Assuming that $q(\cdot;\ )$ is the density associated to $Q(\cdot;\ )$, $Y_1,...,Y_{d-1}$ have density

$$q[F_1(y_1;\ ),...,F_1(y_{d-1};\ );\ ] \prod_{j=1}^{d-1} f_1(y_j;\ ), \qquad (6)$$

where $f_1(y_j;\ )$ is the marginal density of $Y_j$ derived from $F_1(y_j;\ )$.

We then consider a discrete ordinal response. Without loss of generality, we assume that the response is distributed across $K$ consecutive categories $k=1,...,K$ at all times. The joint probability mass function (pmf) $\Pr(Y_1 = y_1,...,Y_{d-1} = y_{d-1};\ \gamma)$ is (Joe, 1997, p. 237)

$$\sum_{l_1=y_1-1}^{y_1} ... \sum_{l_{d-1}=y_{d-1}-1}^{y_{d-1}} (-1)^{\sum_{j=1}^{d-1}(y_j-l_j)} Q[F_1(l_1;\ ),...,F_1(l_{d-1},\ );\ ]. \qquad (7)$$

For simplicity, (7) will also be termed density and referred as $f_{d-1}(y_1,...,y_{d-1};\ \gamma)$.

The conditional density $g_1(y_d/y_1,...,y_{d-1};\ \gamma)$ of the unobserved response $Y_d$ at drop-out does not need to be from the same family as the density $f_{d-1}(y_1,...,y_{d-1};\ \gamma)$ of observed responses. The choice of $g_1$ is merely arbitrary and its adjustment to the non-observed data cannot be verified. As raised by many discussants of Diggle and Kenward (1994), MNAR models rest upon untestable hypotheses and should be interpreted with caution. When MNAR models are considered, sensitivity analyzes are recommended to evaluate the robustness of conclusions to changes in the distributional assumptions. Such analyzes are however not in the scope of this paper.

In our approach, we assume that the marginal distribution $F_1(y_d; )$ is of the same family as distributions at previous times. The conditional distribution $g_1(\cdot; _Y)$ is then computed in the usual way, as the ratio of a joint density $h_d(y_1, ..., y_d)$ of the $d$ responses over $f_{d-1}(y_1, ..., y_{d-1})$. The joint density $h_d(\cdot)$ is constructed in a similar way as $f_{d-1}(\cdot)$, from equations (6) and (7) in the continuous and the ordinal cases respectively.

## 5. Dependence between Drop-out and Response

We specify a marginal Bernoulli model indexed by a vector $_D$ for the drop-out status $O_j$ given that $O_{j-1} = 1$ for $j = 1, ..., d$. The joint distribution of $(O_j/O_{j-1} = 1)$ and of any prior response $Y_l$ ($l \leq j$) is modelled using a copula $K(\cdot; )$ with parameter . The bivariate distribution $F_2(y_l, o_j/O_{j-1} = 1; )$ is defined as

$$K[F_1(y_l; ), \Pr(O_j \leq o_j/O_{j-1} = 1; _D); ]. \tag{8}$$

The conditional distribution of $O_j$ given $y_l$ ($l \leq j$) is then computed from the (discrete) derivation of Equation (8) with respect to $F_1(y_l; )$.

With continuous responses, $\Pr(O_j \leq o_j/O_{j-1} = 1, Y_l = y_l; )$ is (Joe, 1997, p. 245)

$$\left. \frac{\partial K[x, \Pr(O_j \leq o_j/O_{j-1} = 1; _D); ]}{\partial x} \right|_{x = F_1(y_l; )}. \tag{9}$$

With a discrete ordinal response (on $k = 1, ..., K$), it is

$$\frac{\sum_{k_l = y_l - 1}^{y_l} (-1)^{(y_l - k_l)} K[F_1(k_l; ), \Pr(O_j \leq o_j/O_{j-1} = 1; _D); ]}{f_1(y_l; )}.$$

## 6. Some Copula Distributions and Dependence Properties

In this section, we review the copula distributions that will be used in the illustrations. Other parametric families of copulas can be found in Joe (1997,

ch. 5). As defined in (1), copulas are direct functions of the marginal distributions $F_1(z_j)$. They characterize the type and the strength of dependence between margins. When all the $Z_j$ are mutually independent, the joint distribution is the product of the marginal distributions. A particular copula, called the *product*, is used to characterize independence. It is

$$(u_1, \ldots, u_J) = u_1 \ldots u_J. \tag{10}$$

Another copula is created from the multivariate standard normal distribution $_J(\cdot; R)$, with correlation matrix $R$. It is given as

$$C\ (u_1, \ldots, u_J; R) = \ _J[\ _1^{-1}(u_1), \ldots, \ _1^{-1}(u_J); R], \tag{11}$$

where $_1(\cdot)$ is the univariate standard normal distribution.

The dependence between random variables $U_j$ is modelled in (11) using the correlation matrix $R$. Under independence, $R$ is the identity and $C\ (\cdot; R)$ reduces to $\ (\cdot)$. When all marginal distributions $F_1(z_j)$ are normal, (11) generates a multivariate normal distribution for the $Z_j$. The matrix $R$ is then also the correlation matrix for the $Z_j$. When the $Z_j$ are not all normally distributed, $R$ still quantifies dependence but it is not a formal correlation matrix anymore. With continuous responses, the $(i, j)^{\text{th}}$ entry of $R$ can be written in the form $2\sin(\pi r_{ij}/6)$, where $r_{ij}$ is Spearman's rank correlation between $Z_i$ and $Z_j$. Standard dependence structures such as exchangeable or autoregressive can be specified using $R$ (see Lambert and Vandenhende, 2002) . In the heart rate example, the copula $C\ (\cdot; R)$ with $R$ indexed by the vector will be used to model the dependence between responses.

Another copula considered in the examples is the family of Frank (1979) with one parameter . Statistical properties of Frank's copula are given in

10

Nelsen (1986) and Genest (1987). It is equal to

$$-\frac{1}{\alpha} \log\left[1 + \frac{\prod_{j=1}^{J}(e^{-\alpha u_j} - 1)}{(e^{-\alpha} - 1)^{J-1}}\right]. \tag{12}$$

Nelsen (1999) discussed various dependence properties of copulas. Random variables $U_1, ..., U_J$ are termed *positively* dependent when large (or small) values are more likely to occur simultaneously than if they were independent. Then, we have

$$C(u_1, ..., u_J) \geq \prod(u_1, ..., u_J), \quad u_j \in [0, 1].$$

Negative dependence occurs when $C(u_1, ..., u_J) \leq \prod(u_1, ..., u_J), \quad u_j \in [0, 1]$. Both positive and negative associations can be modelled using $\alpha$ in the copula of Frank. Positive dependence occurs when $\alpha > 0$ and negative dependence occurs when $\alpha < 0$. The case $\alpha = 0$ does not produce a copula distribution but it is a limiting case for independence ($C = \prod$). In the bivariate case ($J = 2$), the whole range of dependence down to complete negative association can be modelled. When $J > 2$, however, the set of parameter values yielding negative dependence shrinks as $J \to \infty$, so that only parameter values corresponding to positive degrees of dependence are admissible for all dimensions (Kimberling, 1974). A bivariate Frank's family parameterized by the vector $\psi$ will be used to model the dependence between drop-out and responses in the two examples. The same family with parameters $\phi$ will be used for the dependence between repeated measures in the migraine example.

The indicator of drop-out $O_j$ has been ordered in such a way that a low value (0) indicates drop-out and a large value (1) indicates observation. Positive dependence in Equation (8) implies that dropout (i.e., a small value

11

for $O_j$) is more likely to happen with *small* responses. Negative dependence implies the opposite; that subjects tend to drop with large responses. This parameterization requires some care in the interpretation of the dependence parameters   in the illustrations.

   Copula parameters   and   quantify the dependence between marginal distributions. When all random variables are continuous, these parameters are directly related to rank-based dependence measures such as Kendall's tau or Spearman's rho. This relationship is presented and discussed in Nelsen (1999, ch. 5). With non-continuous random variables, a value for the copula parameter can be associated with several estimates of rank-based dependence measures. As shown by Vandenhende and Lambert (2000) for bivariate ordinal responses, the value of the rank-based statistic also depends on the marginal distributions. When marginal distributions are fixed and when considering ordered copulas, they observe a monotonic relationship between the copula parameter and Kendall's tau-b. This monotonicity supports the interpretation of copula parameters as dependence measures, though an exact quantification of tau-b would also require information from the marginal models.

## 7.   Illustrations

A continuous response (heart rate) is modelled in the first example and a discrete ordinal variable (pain severity) is analyzed in the second example. All computation were performed using the *Mathematica$^{TM}$* package. The statistical add-ons *ContinuousDistributions* and *MultinormalDistribution* were used to access standard univariate and the multivariate normal distributions, respectively. Frank's copula and the discrete distributions were hard-coded.

Maximum likelihood estimates were obtained using the *FindMinimum* built-in function. The observed information matrix was derived using the *D* function and standard errors were calculated using the *Inverse* function. Programs were rather short (2 pages) and fast to run (a few seconds).

## 7.1 *Heart rate in Toxicology Study*

An acute toxicology study was performed in male rats. Five doses of an investigational drug were randomly administered to groups of 6 animals each. Heart rate (HR) was recorded in beats per minute at several times relative to dosing. Results are summarized in Table 1. Mean baseline HR were comparable between groups. Treatment tends to decrease mean HR but the e ect is not dose-related. A few animals died during the study, predominantly at the high dose. The two animals dying within 30 minutes had baseline HR slightly above the average. Thereafter, HR recorded in animals prior to death were below or close to the group's geometric mean at time of measurement. This suggests that a time-varying dependence between drop-out and HR is possible. The following models were fitted to the data.

[Table 1 about here.]

Several marginal distributions (Cauchy, Gamma, Log-normal) were tested and the Log-normal distribution was finally retained for HR because of its fit. A log-linear regression on $\log(t_j)$ was applied to each group separately. Due to the similarity of mean baseline values between groups, a single intercept was estimated across groups. The model was parameterized as $E[\log(Y_j)] = \quad + \quad_{group} \log(t_j)$, where is the intercept and $_{group}$ are dose-specific slopes. The dependence between repeated responses was modelled

13

using the copula defined in (11) for $Q(\cdot;\ )$, with indexing the correlation matrix $R$. A first-order autoregressive - AR(1) - dependence structure was chosen so that the $(j, l)^{th}$ element of $R$ was equal to $\ ^{|j-l|}$ $(0 \leq\ \leq 1)$. Note that in the case of Log-normal margins, the copula (11) approach is identical to a traditional AR(1) model on log-responses with normal innovations. The drop-out hazard was assumed constant over time and modelled using a logistic regression. The logit of $p_j = \Pr(O_j = 0 / O_{j-1} = 1)$ was linearly regressed on the administered dose as $\mathrm{logit}(p_j) = \ _{Dint} + \ _{Ddose}\mathrm{dose}$. Two models for the dependence between drop-out probability and response were developed. A MCAR process was modelled using the product copula in (10). A MAR process was modelled assuming that drop-out probability was conditionally independent of the historical data, given the last prior observation. A bivariate Frank's copula as in (12) was used for $K(\cdot;\ _j)$ to relate $\Pr(O_j \leq o_j / O_{j-1} = 1)$ and $F_1(y_{j-1})$. A distinct copula parameter $_j$ was first fitted for each time. The graphical representation of the estimated $_j$ (with 95% confidence intervals) over time suggested that a linear relationship $_j = \ ^0 + \ ^1 \log(t_{j-1})$ could be considered. The baseline time $(t_1)$ was taken as 1 minute instead of zero for identifiability of the above model. This yields a dependence equal to $\ ^0$ between response and drop-out when occurring prior to 30 min $(j = 2)$ and, then a monotonic change in the dependence over time.

The maximum likelihood estimate (MLE) $\ ^0$ of the MAR model (-2 log lik=-310.1; 12 parameters) is equal to -6.37 (s.e.=5.13). A reduced MAR model with $\ ^0$ fixed to zero was also fitted. This model implies independence between drop-out and responses for deaths prior to 30 min. Parameter esti-

mates from the MCAR and the reduced MAR models are presented in Table (2). The reduced MAR model (-2 loglik=-306.6; 11 parameters) is selected as it performs slightly better than the initial MAR and than the MCAR (-2log lik=-302.2; 10 parameters) models. Estimates and standard errors of parameters common to the three models are comparable (result not shown for the initial MAR model). Compared to placebo, the drug decreases HR but a reversed trend is observed with dose. The decrease is more important at low than at high doses. The drop-out hazard (death rate) increases significantly with dose. Predicted probabilities of death are 2.3, 2.5, 3.0, 5.3 and 29.8 % in the 0, 3, 10, 30 and 100 mg/kg groups respectively. A significant positive estimate of $^1$ is found in the reduced MAR model. As discussed previously, this indicates a negative dependence between HR and death after 30 minutes post-dosing. From that point on, animals with a low HR are more likely to die than those with a high rate and the trend inflates with time.

[Table 2 about here.]

As illustrated in this example, our modelling strategy allows for marginal analyzes of the repeated response and drop-out rates, while controlling for their dependence. The selected MAR model does not permit separability between parameter sets $\{\ ,\ _{dose},\ ^2,\ \}$ and $\{\ _{Dint},\ _{Ddose},\ ^1\}$. Therefore, different MLEs and standard errors for $\{\ ,\ _{dose},\ ^2,\ \}$ are obtained in the MCAR and MAR columns of Table (2). The difference is reasonably small and does not alter conclusions. Parameter sets are also moderately correlated (result not shown) in this example. A maximum correlation of 0.30 is found between estimates of $_{100}$ and $_{Ddose}$.

15

A drop-out probability model which is conditional on the response is also readily available from Equation (9). Such models are useful for the interpretation of the drop-out process, in relation to the response variable. Some conditional drop-out profiles (from the reduced MAR model) are presented for the 100 mg/kg group in Figure 1. At 30 min, the drop-out probability (29.8%) is identical whatever the response. This is the marginal estimate. Thereafter, the drop-out probability increases with decreasing HR. Our model assumes a monotonic change of the dependence parameter $\alpha_j$ over time. However, when considering a particular HR level (e.g., 375 beats/min), we do not observe a monotonic change in conditional probabilities anymore. The drop-out probability peaks at 1h and decreases thereafter. The marginal distribution of HR is also changing over time (decreased mean, constant CV). A HR as low as 375 beats/min is rather atypical in the 100mg/kg group at 1h post-dose but this value is closer to the predicted mean at subsequent times. This explains the reversed trend in conditional drop-out probability after 1h at that level.

[Figure 1 about here.]

Selection models where $f^{D/Y}$ extends the initial logistic drop-out model were also fitted. Conditional models had additional terms for the response at $j - 1$ (-2logL=-304.0; 11 parameters), the change in response between $j - 2$ and $j - 1$ (-2logL=-303.2; 11 parameters), or both the response and the change (-2logL=-304.2; 12 parameters). The reduced MAR copula model provides a better fit to the data than any of these selection models.

7.2 *Clinical Trial on Migraine*

The second example is a clinical trial on acute migraine. Thirty nine patients with a moderate to severe migraine were randomized into 3 treatment

16

groups (0, 5 or 20 mg of LY334370 i.v.). Pain severity (PS) was measured on a four-grade scale (no, mild, moderate, severe) on 8 occasions after dosing. Subjects were allowed to take a rescue medication from 2 hours post-dose but, thereafter, PS was not collected anymore. The cumulative incidence of PS and drop-outs is displayed in the upper panel of Figure 2. There is a progression towards lower severity with time and the PS decrease is more important in the two treated groups than under placebo. The rate of drop-outs is also larger when not treated. A copula-based analysis of these PS data ignoring drop-outs was made by Vandenhende and Lambert (2000). They tested several copula families to model an AR(1) process. Marginal parameter estimates and predicted rank-based dependence measures were not much influenced by the considered copulas. Therefore, we shall use their final model considering Frank's family for the response. Then, we shall model the drop-out hazard and the dependence between drop-out and PS.

[Figure 2 about here.]

Pain severity is an ordinal response with four possible values ($K = 4$). No pain is coded as $k = 1$, mild pain as 2, moderate as 3 and severe as 4. A cumulative regression model with complementary log-log link was selected. The model included 3 threshold parameters $\theta_k$ for the levels $k = 1, 2, 3$, two parameters $\beta_{5-0}$ and $\beta_{20-0}$ for the contrasts between the treated groups (5 and 20 mg respectively) and placebo and a log-time effect $\beta_{time}$. The dependence between successive responses was modelled based on a bivariate Frank's copula with parameter $\alpha$ and the joint distribution of responses was constructed assuming a first order Markov process. This induces an AR(1)

17

dependence structure, with dependence decreasing with lag time. Drop-out hazard was modelled using a time-independent logistic regression including a different parameter for placebo ($_0$) and for drug ($_{5,20}$). Dependence models for MCAR, MAR and MNAR drop-out processes were considered. The MAR and MNAR models were created from the bivariate Frank's copula with parameter  under a conditional independence hypothesis. The MAR model relates $\Pr(O_j \leq o_j | O_{j-1} = 1)$ to $F_1(y_{j-1})$ as in the previous example. The MNAR model relates it to $F_1(y_j)$, the distribution of the possibly un-observed response at $t_j$. All marginal distributions of the response $F_1(\cdot)$ are multinomial and the distribution at drop-out is extrapolated from the distributions of previously observed data using the same AR(1) model as above. Parameters were jointly estimated. MLEs are in Table 3.

[Table 3 about here.]

MLEs and standard errors of parameters common to the three models are comparable. Compared to placebo, LY334370 reduces PS and the effect is comparable between doses. A strong positive dependence is found between successive responses (Kendall's tau-b $\approx$ 0.75). The drop-out rate is larger under placebo (14%) than under active therapy (4%). These estimates are from the MAR model (-2 loglik=427.1; 10 parameters), which performs better than the MCAR (-2 loglik=443.6; 9 parameters) and the MNAR (-2 loglik=429.9; 10 parameters) models. A positive dependence between drop-out and response is detected (Kendall's tau-b $\approx$ 50%), implying than patients are more likely to take a rescue medication when their PS is large. The predicted cumulative distributions of PS and drop-out are in the lower panel

of Figure 2. Predicted and empirical (see the upper panel) distributions are comparable. This provides evidence for the good adjustment of the MAR model to the data.

## 8. Discussion

Our interest was in the specification of marginal models for repeated responses and drop-out events in longitudinal studies with non-random attrition. The dependence between successive responses and between drop-out and response was modelled using copulas and a full likelihood-based method was proposed.

A series of parametric copula families is available; see Joe (1997, ch. 5). The selection of a candidate is either made by convenience (e.g. when choosing the multivariate normal distribution) or based on the comparison of the adjustment of different models using likelihood ratios (see e.g. Vandenhende and Lambert, 2000). Genest and Rivest (1993) have developed a goodness of fit test and a graphical method to select copula models in comparison to the empirical dependence found in the data. Their method applies to a subfamily of bivariate copulas named *Archimedean*. Extensions to other families and to the multivariate case would be worth considering.

The issue of sensitivity is often raised when dealing with non-random drop-outs. As illustrated by Kenward (1998), several selection models with a comparable adjustment can yield very different conclusions. This is related to the untestable distributional assumption for the unobserved responses at time of drop-out. In the HR example, we decided not to further investigate MNAR models for this reason. Here, drop-out is synonym of death and model predictions (from data collected in a living state) were in no case satisfactory

to meet the physiological inactivation of the heart at death. The situation in the PS example was somewhat di erent in the sense that PS was still predictable around the intake of a rescue medication. As noted by Kenward (1998), *"the MNAR analysis tells us about inadequacies of the original model rather than the adequacy of the MNAR model"*. In the example, parameters estimates were comparable between models.

Our copula models relate the drop-out process to the *distribution* of response before (MAR) or at (MNAR) time of drop-out. This contrasts with standard specifications of selection models, in which the drop-out hazard is modelled in relation to the *actual level* of the response (using e.g. logistic regressions). When the distribution of response is evolving over time, the two specifications provide distinct natural interpretations for the conditional drop-out process. Let us illustrate that point in the migraine trial example. During a classical migraine crisis, PS tends to decrease over time and usually terminates within 48 hours. It is also natural to assume that the probability of drop-out (RM intake) increases with pain intensity. However, the risk of drop-out for a patient with a mild pain severity will probably be larger at 24 hour than at 2 hour post-dosing. This is because mild pain is not considered as unacceptable (abnormal value) at 2 hours, whereas it is rather atypical and more worrying when still present 24-hour after the crisis onset. Such a time by response interaction is often neglected but would need to be modelled in selection models. With copulas, no such modelling is needed. The conditional drop-out probability model is automatically adjusted for the distribution of responses at the time of drop-out. This behavior is further illustrated in Figure 1 for the heart rate example.

20

## Acknowledgements

## References

Diggle, P. J. and Kenward, M. G. (1994). Informative drop-out in longitudinal data analysis. *Applied Statistics* 43, 49–93.

Frank, M. J. (1979). On the simultaneous associativity of F(x,y) and x+y-F(x,y). *Aequationes Mathematicae* 19, 194–226.

Genest, C. (1987). Frank's family of bivariate distributions. *Biometrika* 74, 549–555.

Genest, C., Ghoudi, K. and Rivest, L.-P. (1995). A semiparametric estimation procedure of dependence parameters in multivariate families of distributions. *Biometrika* 82, 543–552.

Genest, C. and Rivest, L.-P. (1993). Statistical inference procedures for bivariate Archimedean copulas. *Journal of the American Statistical Association* 88, 1034–1043.

Joe, H. (1997). *Multivariate Models and Dependence Concepts*. Chapman & Hall.

Kenward, M. G. (1998). Selection models for repeated measurements with non-random dropout: An illustration of sensitivity. *Statistics in Medicine* 17, 2723–2732.

Kimberling, C. H. (1974). A probabilistic interpretation of complete mono-
tonicity. *Aequationes Mathematicae* 10, 152–164.

Lambert, P. and Vandenhende, F. (2002). A copula based model for multi-
variate non normal longitudinal data: analysis of a dose titration safety
study on a new antidepressant. *Statistics in Medicine* (in press).

Little, R. J. A. (1993). Pattern-mixture models for multivariate incomplete
data. *Journal of the American Statistical Association* 88, 125–134.

Little, R. J. A. (1994). A class of pattern-mixture models for normal incom-
plete data. *Biometrika* 81, 471–483.

Little, R. J. A. and Rubin, R. B. (1987). *Statistical Analysis with Missing
Data.* Wiley.

Meester, S. G. and MacKay, J. (1994). A parametric model for clustered
correlated categorical data. *Biometrics* 50, 954–963.

Molenberghs, G., Kenward, M. G. and Lesa re, E. (1997). The analysis
of longitudinal ordinal data with nonrandom drop-out. *Biometrika* 84,
33–44.

Molenberghs, G., Michiels, B. and Kenward, M. G. (1998). Pseudo-likelihood
for combined selection and pattern-mixture models for incomplete data.
*Biometrical Journal* 40, 557–572.

Molenberghs, G., Michiels, B., Kenward, M. G. and Diggle, P. J. (1998).
Monotone missing data and pattern-mixture models. *Statistica Neer-
landica* 52, 153–161.

Nelsen, R. B. (1986). Properties of a one-parameter family of distributions
with specified marginals. *Communications in Statistics, Theory and Meth-
ods* 15, 3277–3285.

Nelsen, R. B. (1999). *An Introduction to Copulas*. Springer-Verlag.

Oakes, D. (1994). Multivariate survival distributions. *Nonparametric Statistics* 3, 343–354.

Rubin, D. B. (1976). Inference and missing data. *Biometrika* 63, 581–592.

Shih, J. H. and Louis, T. A. (1995). Inferences on the association parameter in copula models for bivariate survival data. *Biometrics* 51, 1384–1399.

Sklar, M. (1959). Fonctions de répartition à n dimensions et leurs marges. *Publication Institut Statistiques Université Paris* 8, 229–231.

Trégouèt, D.-A., Ducimetière, P., Bocquet, V., Visvikis, S., Soubrier, F. and Tiret, L. (1999). A parametric copula model for analysis of familial binary data. *American Journal of Human Genetics* 64, 886–893.

Vandenhende, F. and Lambert, P. (2000). Modeling repeated ordered categorical data using copulas. *Discussion Paper* 00–25. Institut de Statistique, Université Catholique de Louvain, Louvain-la-Neuve, Belgium *(ftp://www.stat.ucl.ac.be/pub/papers/dp/dp00/dp0025.ps)*.

Wang, W. and Wells, M. T. (2000). Model selection and semiparametric inference for bivariate failure-time data. *Journal of the American Statistical Association* 95, 62–76.
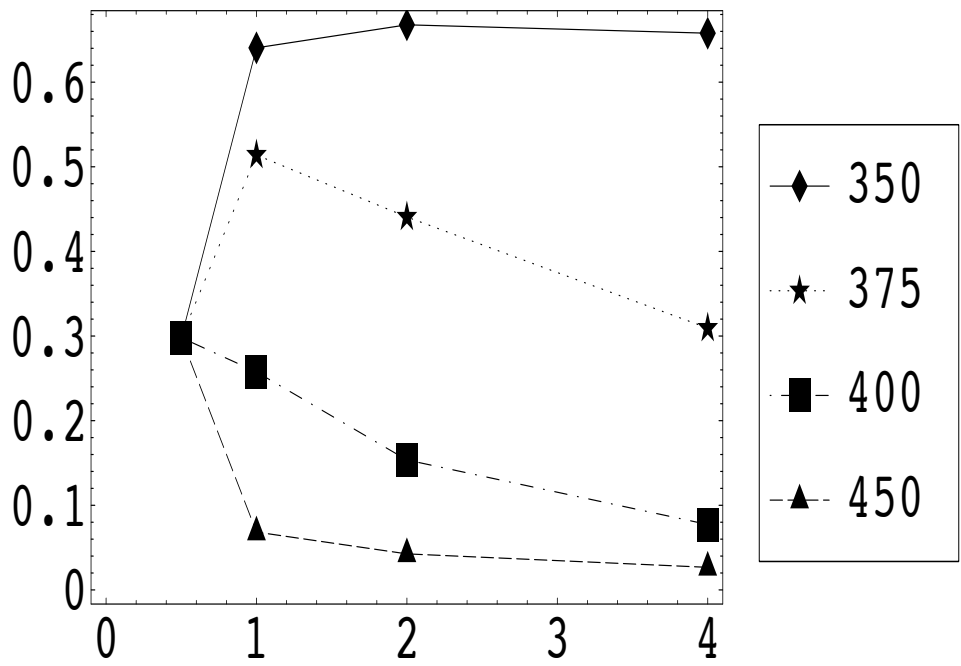
Figure 1. Predicted conditional drop-out hazard profiles (v.s. time (h)) in the 100 mg/kg group from the MAR model of Table (2). The marginal hazard (29.8 %) is plotted at 0.5h, when drop-out is assumed independent of heart rate. Lines are drawn over time for constant heart rates of 350, 375, 400 and 450 beats/min.
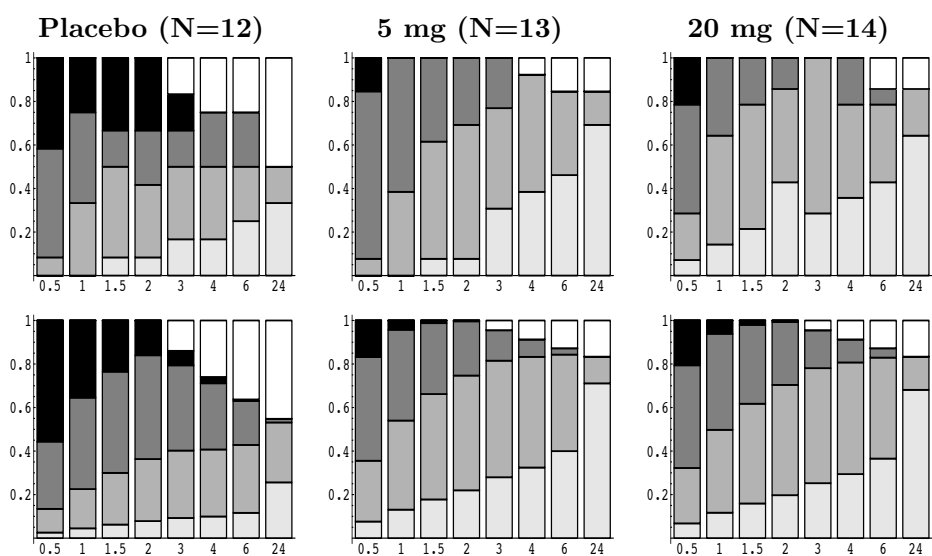
Figure 2. Empirical (upper panel) and predicted (lower panel) cumulative distributions of pain severity: no (10% gray), mild (30% gray), moderate (50% gray), severe (black) and drop-out (white) over time in the three treatment groups. Predictions are obtained from estimates of the MAR model in Table (3).

Table 1

*Summary statistics on heart rate (beats/min) in male rats after intake of an investigational drug at dose levels of 0, 3, 10, 30 and 100 mg/kg. Six animals were randomized to each group but some animals died during the study. Heart rates taken before animals' death are in parentheses.*

| Dose | (mg/kg) | Time (min) after dosing | | | | |
|---|---|---|---|---|---|---|
| | | 0 | 30 | 60 | 120 | 240 |
| 0 | N | 6 | 6 | 6 | 6 | 6 |
| | Mean[1] | 461 | 495 | 481 | 524 | 513 |
| | CV(%) | 11 | 12 | 8 | 7 | 4 |
| | | | | | | |
| 3 | N | 6 | 6 | 6 | 6 | 5 |
| | | | | | (330) | |
| | Mean | 458 | 374 | 355 | 335 | 331 |
| | CV(%) | 11 | 9 | 4 | 4 | 7 |
| | | | | | | |
| 10 | N | 6 | 6 | 6 | 5 | 5 |
| | | | | (370) | | |
| | Mean | 454 | 369 | 361 | 352 | 338 |
| | CV(%) | 6 | 7 | 6 | 4 | 1 |
| | | | | | | |
| 30 | N | 6 | 6 | 5 | 5 | 5 |
| | | | (330) | | | |
| | Mean | 479 | 373 | 382 | 382 | 356 |
| | CV(%) | 8 | 10 | 5 | 2 | 4 |
| | | | | | | |
| 100 | N | 6 | 4 | 2 | 2 | 2 |
| | | (490,500) | (390,380) | | | |
| | Mean | 457 | 397 | 440 | 395 | 380 |
| | CV(%) | 8 | 7 | 6 | 2 | 0 |

1: Geometric means are reported.

**Table 2**

*Maximum likelihood estimates (s.e.) for the joint repeated-measures and survival analyzes of the heart rate data.*

| Parameter | MCAR Model (-2 loglik=-302.2) | MAR Model (-2 loglik=-306.7) |
|---|---|---|
| Marginal model on heart rate | | |
| $\alpha$ | 6.14 (0.013) | 6.14 (0.013) |
| $\alpha_0$ | 0.0194 (0.0047) | 0.0197 (0.0047) |
| $\alpha_3$ | -0.0629 (0.0048) | -0.0616 (0.0047) |
| $\alpha_{10}$ | -0.0588 (0.0049) | -0.0575 (0.0049) |
| $\alpha_{30}$ | -0.0522 (0.0052) | -0.0524 (0.0052) |
| $\alpha_{100}$ | -0.0341 (0.0066) | -0.0353 (0.0062) |
| $\sigma^2$ | 0.0049 (0.0007) | 0.0049 (0.0007) |
| | | |
| Dependence between successive heart rates | | |
| $\rho$ | 0.442 (0.091) | 0.450 (0.091) |
| | | |
| Marginal model of drop-out hazard | | |
| $\beta_{Dint}$ | -3.75 (0.68) | -3.76 (0.68) |
| $\beta_{Ddose}$ | 0.0283 (0.010) | 0.0290 (0.009) |
| | | |
| Dependence between drop-out and heart rate | | |
| $\theta^1$ | | 1.01 (0.56) |

## Table 3
*Maximum likelihood estimates (s.e.) for the joint repeated-measures and survival analyzes of the migraine data.*

| Parameter | MCAR model (-2 loglik=443.6) | MAR model (-2 loglik=427.1) | MNAR model (-2 loglik=429.9) |
|---|---|---|---|
| Marginal model on pain severity score | | | |
| $\alpha_1$ | -3.13 (0.40) | -3.08 (0.40) | -3.09 (0.40) |
| $\alpha_2$ | -1.44 (0.29) | -1.37 (0.29) | -1.46 (0.29) |
| $\alpha_3$ | -0.0066 (0.23) | 0.036 (0.23) | -0.07 (0.24) |
| $\alpha_{time}$ | 0.81 (0.10) | 0.82 (0.10) | 0.75 (0.09) |
| $\alpha_{5-0}$ | 1.19 (0.31) | 1.11 (0.31) | 1.22 (0.31) |
| $\alpha_{20-0}$ | 1.10 (0.32) | 0.99 (0.32) | 1.05 (0.33) |
| | | | |
| Dependence between successive scores | | | |
| $\rho$ | 13.9 (2.14) | 13.8 (2.12) | 14.5 (2.22) |
| | | | |
| Marginal model on rescue medication (RM) intake | | | |
| $\beta_0$ | -1.74 (0.44) | -1.82 (0.48) | -1.79 (0.52) |
| $\beta_{5,20}$ | -3.19 (0.51) | -3.07 (0.52) | -3.16 (0.53) |
| | | | |
| Dependence between RM intake and pain | | | |
| $\theta$ | | -6.88 (2.76) | -11.9 (8.96) |