

Longitudinal Data Analysis

ANALYSIS with DROP-OUT

Analysis of Longitudinal Data Subject to Drop-out

- ★ One issue in the analysis of longitudinal data that requires attention is the nature of any missing data.
- Missing data can bias results!
- Often we'd prefer to analyze the data with the missing values recovered. However, this isn't always the case (ie. missing due to death).
- ★ There are several statistical approaches that have been proposed for the analysis of longitudinal data subject to drop-out.

Analysis with Missing Data

Outline:

- Example of missing data
- The “Taxonomy” for missing data
- Impact of missing data
- Analysis approaches
- Caveats

Overview: Longitudinal Data with Attrition

- Classification of Missing Data
 - Little & Rubin (1987)
 - Laird (1988)
- Selection Models & Likelihood Analysis
 - Diggle & Kenward (1994)
 - Diggle (1998)

Overview: Longitudinal Data with Attrition

- Pattern Mixture Models
 - Little (1993)
- Inverse Probability Weighted Analysis
 - Robins, Rotnitzky & Zhao (1995)

Overview: Longitudinal Data with Attrition

- Imputation Methods
 - Paik (1997)
- Latent Variable Models
 - Wu and Carroll (1988)
 - Tenhave, Kunselman, Pulkstenis, and Landis (1998)

Examples of Missing Data

- Schizophrenia Treatment Trial
 - ▶ A randomized longitudinal study of haloperidol and risperidone
 - ▶ Primary outcome: PANSS score at 8 weeks
 - ▶ Intermediate outcomes at 1, 2, 4, and 6 weeks
 - ▶ 8 week completion by arm:

Placebo $27/88 = 31\%$

Haloperidol $36/87 = 41\%$

Risperidone (6mg) $52/86 = 60\%$

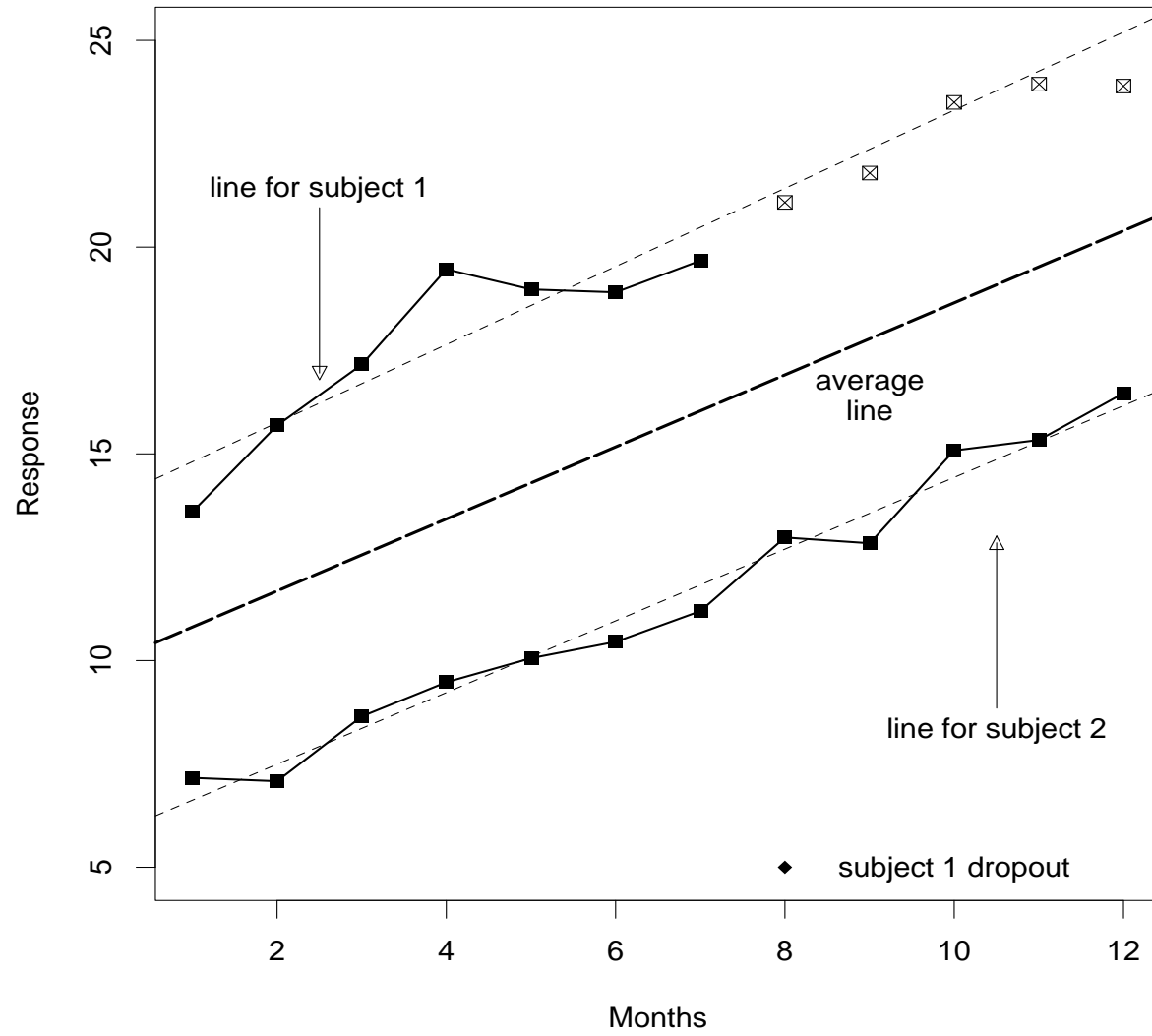
Schizophrenia Treatment Trial

- Reasons for dropout:

Abnormal lab result	4
Adverse experience	26
Inadequate response	183
Inter-current illness	3
Lost to follow-up	3
Uncooperative	25
Withdrew consent	19
Other	7

- This combines the 6 treatment arms

Two Subjects

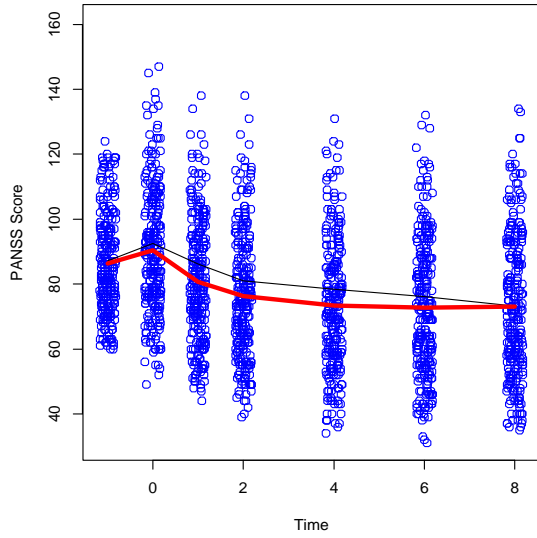


Schizophrenia Trial: Implications of Missing Data

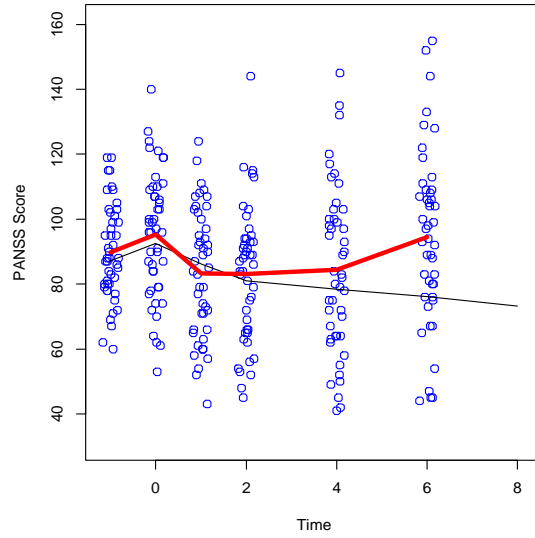
Q: Can we simply use the available cases at 8 weeks to make inference regarding the treatment effect?

Q: Do we have any information that can be used to predict the missing information?

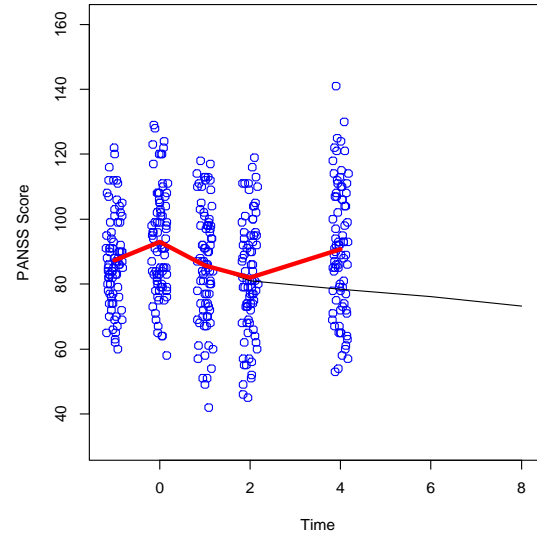
Last Visit = 8



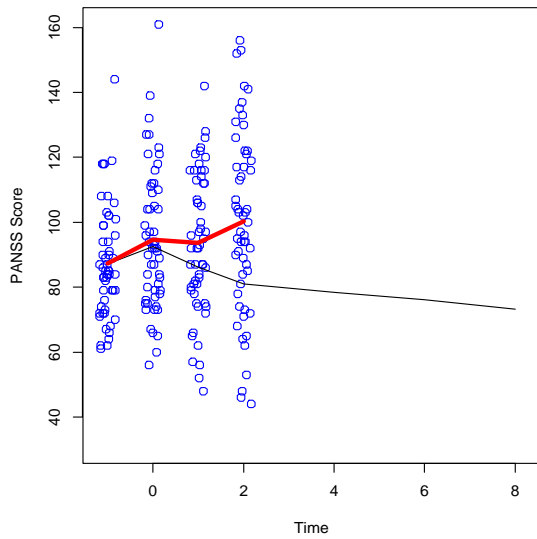
Last Visit = 6



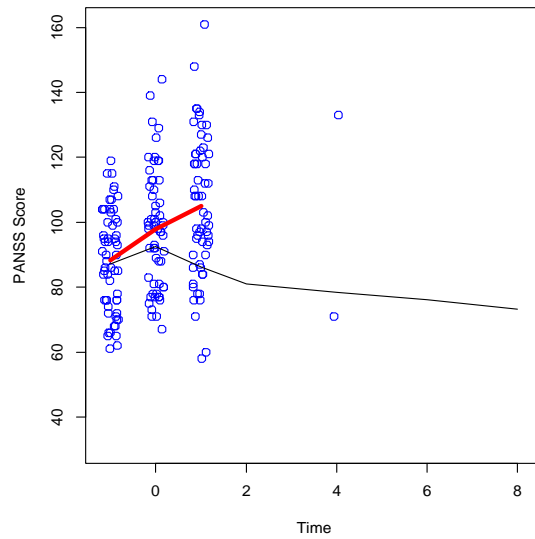
Last Visit = 4



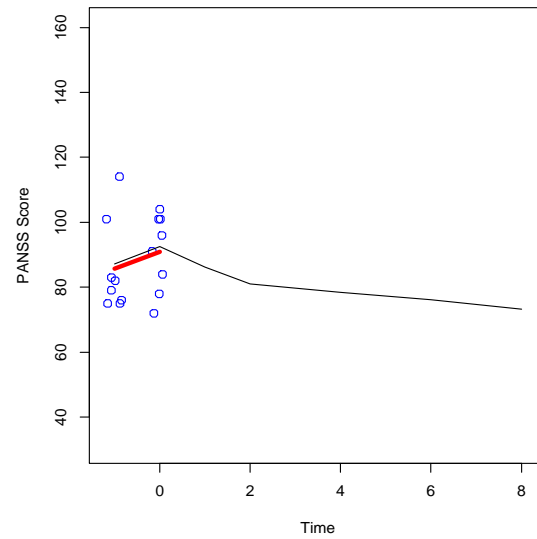
Last Visit = 2



Last Visit = 1



Last Visit = 0



PANSS Data: Dropout prediction

- Regress time-until-dropout on time-varying PANSS score
 - ▷ level = $[Y(t - 1) + Y(t - 2)]/2$
 - ▷ trend = $[Y(t - 1) - Y(t - 2)]/2$

variable	coefficient	s.e.	<i>p</i> -value
level/10	0.141	(0.020)	$p < 0.001$
trend/10	0.499	(0.055)	$p < 0.001$

Some simple approaches

- Last observation carried forward (LOCF)
 - ▷ This was used for publication of schizophrenia trial
 - ▷ Probably helps alleviate bias here (see pattern plot)
 - ▷ There are “better” methods for using past to predict future (Mixed Models; Multiple Imputation; Weighting)
- Complete case analysis
 - ▷ Selection bias
- Derive a composite end-point
 - ▷ Eg. failure = quit or did not improve by 20%
 - ▷ Changes the question
 - ▷ Efficacy, intent-to-treat

Some Plausible Perspectives on PANSS Data

- (1) The study was terrible and can't be used for any reliable inference.
- (2) They should have at least attempted to get an 8 week measurement on each person and then they could do intent-to-treat analysis. (analysis of treatment program)
- (3) Removal of participants based on their observed past outcomes is an example of "missing at random" data and therefore longitudinal analysis using mixed models could provide valid inference.

Taxonomy of Missing Data

Missing Completely at Random (MCAR)

- This assumes that the probability of missing an observation does not depend on any variables. No selection bias.

Missing at Random (MAR)

- This assumes that missing an observation is predicted by variables that you have measured, but not further dependent on variables you have not measured. Possibly fixable selection bias.

Nonignorable (NI)

- This assumes missing an observation is predicted by variables that you have not measured such as the outcome of interest. Not fixable selection bias!!!

Taxonomy of Missing Data: Examples

- **Missing Completely at Random** (MCAR)
 - ▷ Lightning storm destroys some data.
- **Missing at Random** (MAR)
 - ▷ Based on exam (scores) physician removes.
 - ▷ Subjects far away less likely to return.
- **Nonignorable** (NI)
 - ▷ Subjects currently too sick do not return.

Taxonomy of Missing Data (Koepsell)

Missing Completely at Random (MCAR)

- Missingness entirely unrelated to *true data value*.

Missing at Random (MAR)

- Missingness unrelated to true data value *after accounting for other known characteristics of subject*.

Nonignorable (NI)

- Missing values differ systematically from known values, even after accounting for known subject characteristics.

Some Terminology

Complete Data: The scheduled measurements. This is the outcome vector, \mathbf{Y}_i , that would have been recorded if no missing data occurred.

Missing Data Indicators: The binary variables, $\mathbf{R}_i = \text{vec}(R_{ij})$, that indicate whether Y_{ij} was observed, $R_{ij} = 1$, or unobserved, $R_{ij} = 0$.

Observed (Response) Data: \mathbf{Y}_i^O : Y_{ij} such that $R_{ij} = 1$.

Missing Data: \mathbf{Y}_i^M : Y_{ij} such that $R_{ij} = 0$.

Drop-out Time: If missingness is monotone (ie. once $R_{ij} = 0$ then $R_{ik} = 0 \forall k > j$), then we can define the time that the subject drops out of the study as: $D_i = \min_k (R_{ik} = 0)$.

Missing Data Issues

To formulate different missing data mechanisms we use the notation:

$R_{ij} = 1$ if subject i is observed at time j

$R_{ij} = 0$ if subject i is not observed at time j

MCAR Missing completely at random if

$$P(\mathbf{R}_i | \mathbf{Y}_i, \mathbf{X}_i) = P(\mathbf{R}_i | \mathbf{X}_i)$$

This implies that $E(Y_{ij} | R_{ij} = 1, \mathbf{X}_i) = E(Y_{ij} | \mathbf{X}_i)$.

Missing Data Issues

MAR Missing at random if

$$P(\mathbf{R}_i | \mathbf{Y}_i^O, \mathbf{Y}_i^M, \mathbf{X}_i) = P(\mathbf{R}_i | \mathbf{Y}_i^O, \mathbf{X}_i)$$

Here the probability of missing data only depends on the observed values and not the missing values.

Trouble starts here since this implies

$$E(Y_{ij} | R_{ij} = 1, \mathbf{X}_i) \neq E(Y_{ij} | \mathbf{X}_i) \text{ (possibly).}$$

NI Non-ignorable if

$$P(\mathbf{R}_i | \mathbf{Y}_i^O, \mathbf{Y}_i^M, \mathbf{X}_i) \text{ depends on } \mathbf{Y}_i^M$$

Implications of Taxonomy

Missing Completely at Random (MCAR)

- Standard analysis using the available cases is valid.

Missing at Random (MAR)

- Standard analysis of the available cases is potentially biased. However, there are methods that can provide valid analysis, but these require additional (correct) statistical modelling.

Nonignorable (NI)

- Standard analysis of the available cases is likely biased. The bias can not be corrected since missingness depends on unobserved data and therefore can not be empirically modelled. The recommended approach is to conduct sensitivity analyses:

“It’s bad, but how bad could it be?”

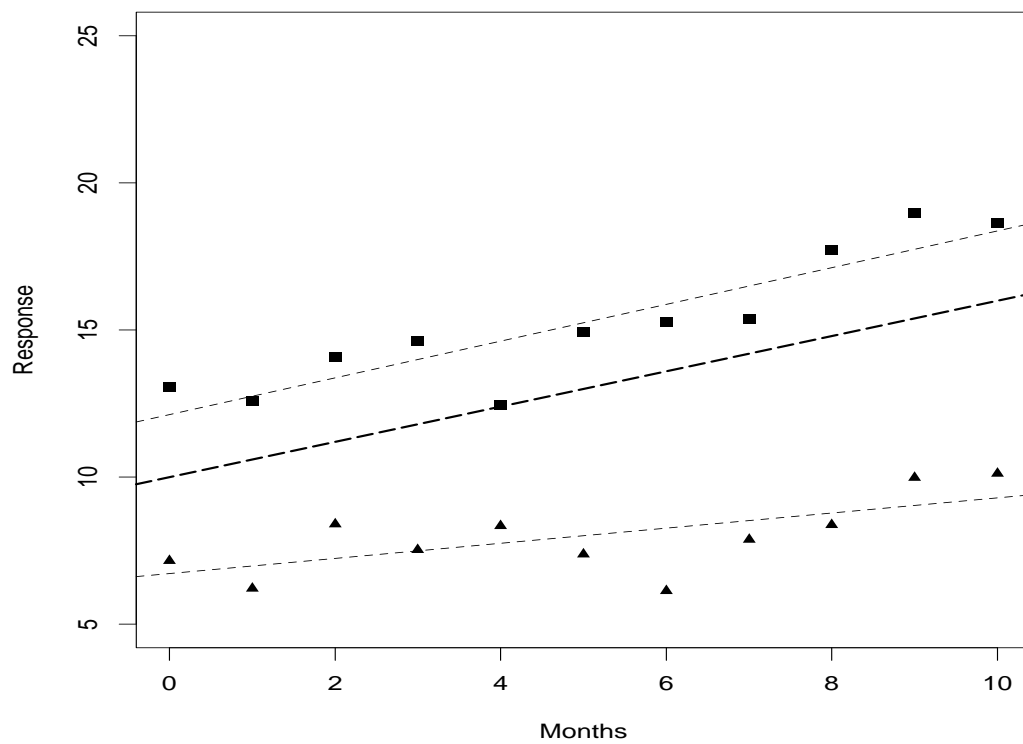
Analysis Approaches when MAR

Maximum likelihood (ML)

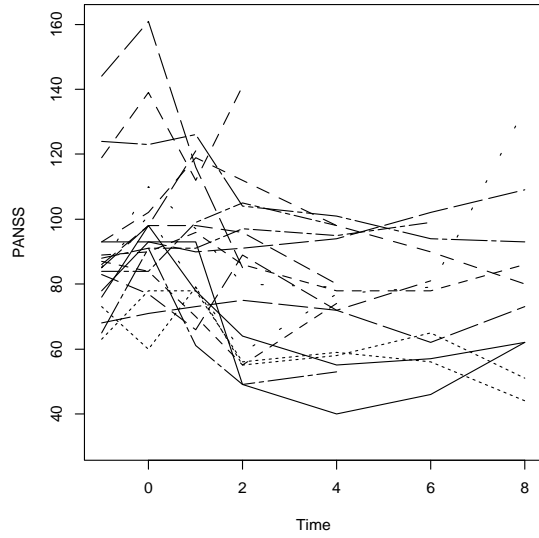
- A model is needed that links the missing outcomes to the factors that predict the missing outcomes.
- The factors that predict the missing outcomes need not be part of the “regression” model of interest (such as intermediate outcomes in a longitudinal study)
- For longitudinal data linear mixed models can help (SAS PROC MIXED)
- This general approach can also be taken for missing covariates.

Linear Mixed Model

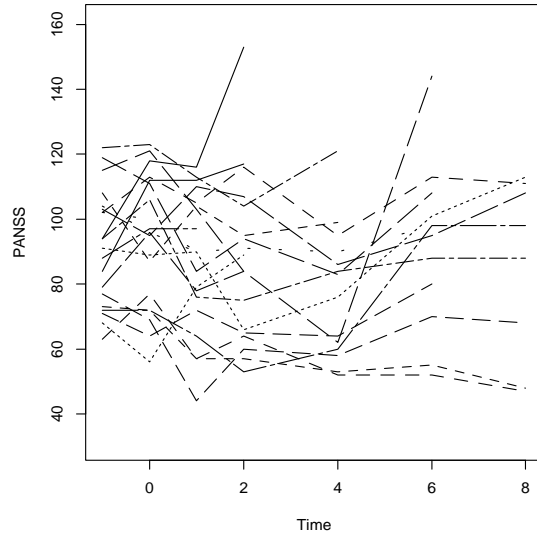
Two Subjects



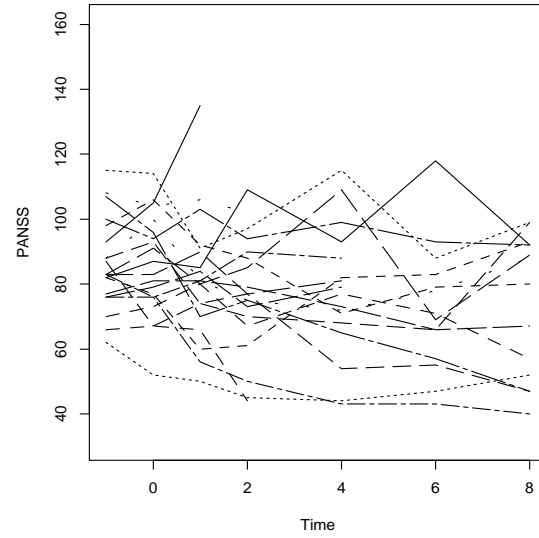
haloperidol



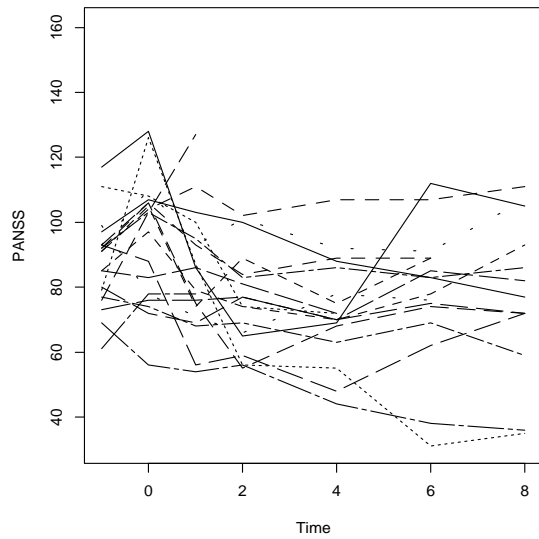
placebo



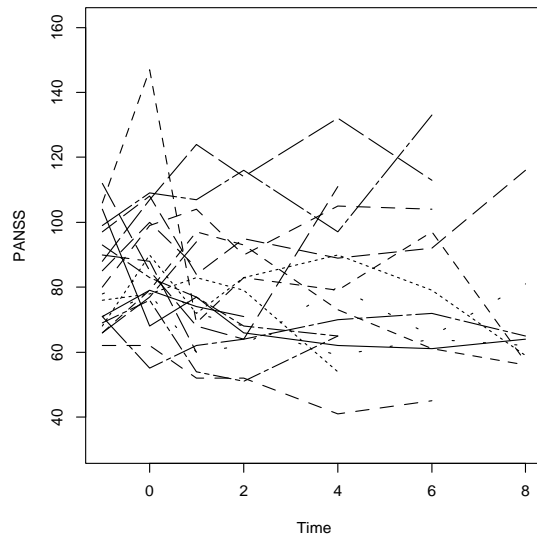
risperidone10



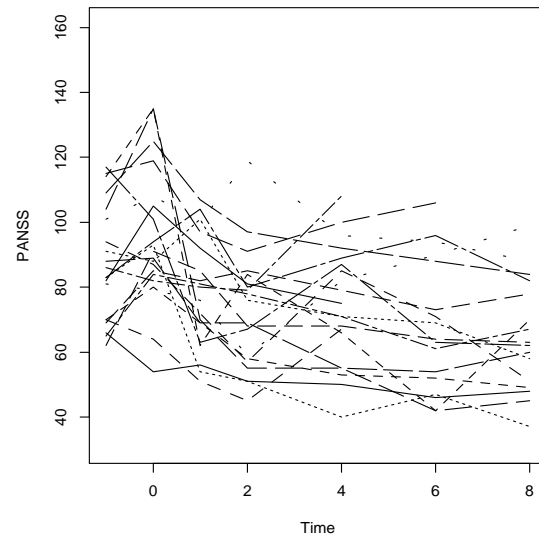
risperidone16



risperidone2



risperidone6



Analysis of PANSS Data

	Final Outcome (8wks)		
	est.	(s.e.)	p-value
Haloperidol is reference group			
Risperidone (6mg)	-7.901	(4.148)	0.0580
Placebo	3.962	(4.844)	0.4142

	Longitudinal Analysis		
	est.	(s.e.)	p-value
Haloperidol is reference group			
Risperidone (6mg)	-13.670	(4.256)	0.0014
Placebo	13.085	(4.586)	0.0045

PANSS Summary

- Analysis using available-case data at 8 weeks yields different conclusions from a longitudinal analysis.

- page 394 suggests dropout predicts the outcome:

$$E[Y(8) \mid \text{completer}] < E[Y(8) \mid \text{drop-out}]$$

- page 395 suggests the outcome predicts drop-out:

$$P[\text{drop-out at } t \mid Y(t-1), Y(t-2)]$$

- Therefore, the data are not **MCAR**.
- **Q**: What is the anticipated impact when analyzing data that are **MAR**?

MAR Illustration: Linear Mixed Models

Response Model $n_i = 12$

$$\mathbf{Y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{b}_i + \mathbf{e}_i$$

$$\mathbf{X}_i = [1, \text{time}, \text{group}, \text{group} \times \text{time}]$$

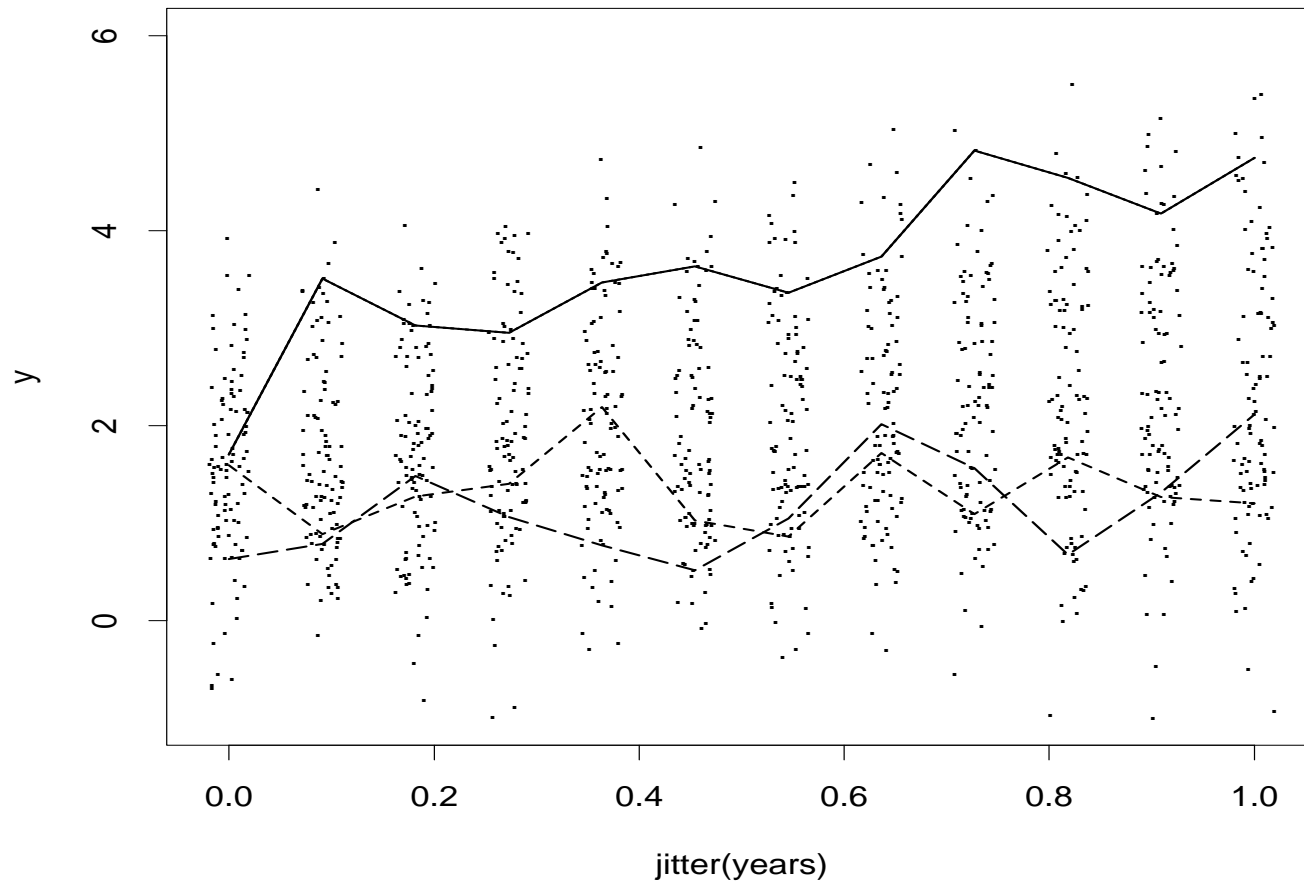
$$\boldsymbol{\beta} = (1.5, 1.0, 0.25, -0.50)$$

$$\mathbf{Z}_i = [1, \text{time}]$$

$$\mathbf{b}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{D}) \quad e_{ij} \sim \mathcal{N}(0, \sigma^2 = 0.2)$$

$$\mathbf{D} = \begin{bmatrix} 0.80 & 0.24 \\ 0.24 & 0.30 \end{bmatrix}$$

Random Lines Data: Complete



MAR Illustration: Linear Mixed Models

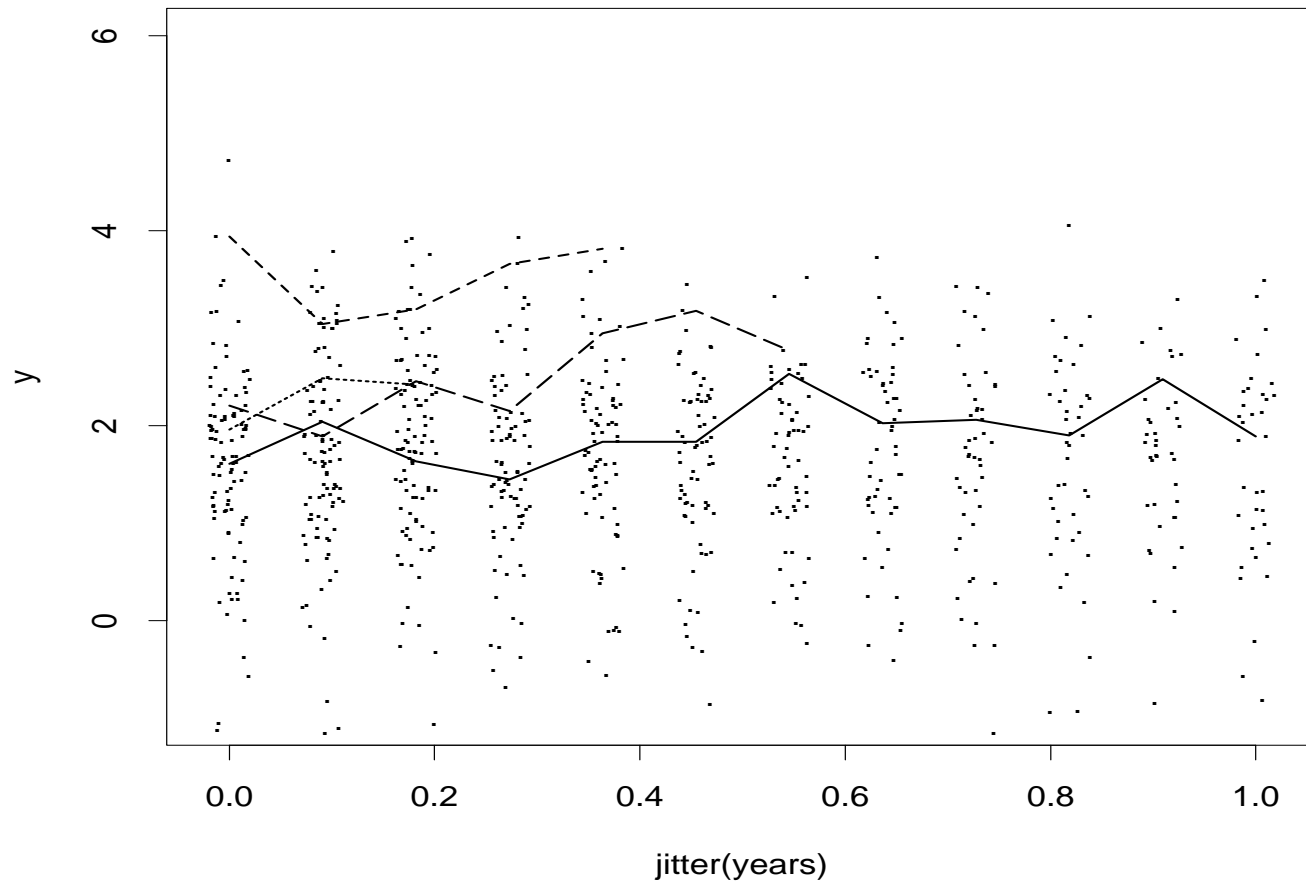
Dropout (Selection) Model

- We use a discrete-time survival model where the hazard for drop-out depends on the previous response values:

$$\begin{aligned} \text{logit}P(R_{ij} = 0 \mid R_{ij-1} = 1, Y_{ik} \ k < j) &= -2.0 \\ &+ 1.0 \cdot (Y_{ij-1} - 2) \\ &+ 0.5 \cdot (Y_{ij-2} - 2) \end{aligned}$$

- $E(n_i) \approx 8$, with 40% of the subjects having complete data.
- N0 = number of subjects with group=0 is 50. N1 (# with group=1) is also 50.

Random Lines Data: MAR



MAR Illustration: Linear Mixed Models

Simulation Study

- A simulation study generated data from the response model and then subjected it to MAR.
- A regression coefficient was estimated using: IEE, GEE-exch, GEE-AR, LMM-(int), and LMM(int+slope).

Mean estimate from 500 simulations

	Method				
	IEE	GEE exch	GEE AR	LMM (int)	LMM (int+slope)
(Int)	1.63	1.55	1.56	1.55	1.52
time	-0.08	0.79	0.40	0.79	1.00
group	0.22	0.24	0.24	0.24	0.24
group.time	-0.53	-0.51	-0.54	-0.51	-0.50

MAR Illustration: Longitudinal Binary Data

Simulation Study

- A similar simulation study generated binary data from a MTM(2) response model and then subjected it to MAR.
- A regression coefficient was estimated using: IEE, GEE-exch, GEE-AR, MTM(1), and MTM(2)
- MTM = “marginalized transition model”. This is likelihood-based estimation for binary longitudinal data (Heagerty 2002) – see Hogan et al. (2004) for an example with use of MTM for data with dropout; see DHLZ (2002) chapter 11 for details.

○ Response Model MTM(2)

$$\text{logit}E[Y_{ij} | \mathbf{X}_i] = -1.5 + 0.5 \cdot \text{time} \\ + 0.5 \cdot \text{group} - 0.25 \cdot \text{group.time}$$

$$\alpha_1 = 2.5$$

$$\alpha_2 = 1.0$$

$$n_i \in [3, 20]$$

$$N0 = N1 = 100$$

- Dropout (Selection) Model

$$\pi_{ij} = P(R_{ij} = 0 \mid R_{ij-1} = 1, Y_{ik} \ k < j)$$

$$\text{logit}(\pi_{ij}) = -3.0 + 1.5 \cdot Y_{ij-1} + 1.0 \cdot Y_{ij-2}$$

- $E(n_i) \approx 11$, with 20% of the subjects having complete data.

Mean estimate from 1000 simulations

	Method				
	IEE	GEE exch	GEE AR	MTM(1) ML	MTM(2) ML
(Int)	-1.56	-1.55	-1.52	-1.52	-1.51
time	0.09	1.26	0.32	0.29	0.51
group	0.51	0.49	0.52	0.52	0.52
group.time	-0.38	-0.12	-0.34	-0.32	-0.26

Example: Complete Data

```
subject 10 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1
subject 11 : 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
subject 12 : 0 0 0 0 0 0 1 1 1 1 1 0 1 1 0 0 0 0 0 0
subject 13 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0
subject 14 : 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 0 1 1 0 0
subject 15 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
subject 16 : 0 0 0 0 0 1 0 0 0 1 1 1 1 0 0 0 0 0 0 0
subject 17 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1
subject 18 : 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0
subject 19 : 0 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1
subject 20 : 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0
subject 21 : 1 1 1 1 1 1 1 0 0 0 0 0 0 0 1 0 1 1 0 0
subject 22 : 0 0 1 0 0 0 0 0 1 0 1 0 0 0 0 0 1 0 0 0
subject 23 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
subject 24 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
subject 25 : 0 1 1 1 1 0 0 0 0 0 1 1 0 0 0 0 0 0 0 1
```

Example: Observed Data

```
subject 10 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1
subject 11 : 1 1 0 0 0 0 0 0 X X X X X X X X X X X
subject 12 : 0 0 0 X X X X X X X X X X X X X X X X
subject 13 : 0 0 0 0 0 0 0 0 0 0 0 X X X X X X X X X
subject 14 : 0 0 0 0 0 0 0 1 1 1 1 1 X X X X X X X X
subject 15 : 0 0 0 0 0 0 0 X X X X X X X X X X X X X
subject 16 : 0 0 0 0 0 1 0 0 0 1 1 1 X X X X X X X X
subject 17 : 0 0 0 0 0 0 0 0 0 0 X X X X X X X X X X
subject 18 : 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 1 X X X X
subject 19 : 0 1 1 1 1 X X X X X X X X X X X X X X X
subject 20 : 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 0
subject 21 : 1 1 1 X X X X X X X X X X X X X X X X X
subject 22 : 0 0 1 X X X X X X X X X X X X X X X X X
subject 23 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
subject 24 : 0 0 0 0 0 0 0 0 0 0 0 0 0 0 X X X X X X
subject 25 : 0 1 1 1 X X X X X X X X X X X X X X X X
```

Analysis Approaches when MAR

Multiple Imputation (MI)

- A model is needed that links the missing variables to the factors that predict the missing outcomes.
- The factors that predict the missing variables need not be part of the “regression” model of interest (such as intermediate outcomes in a longitudinal study)
- Fill in the missing data.
- Fill in and analyze multiple times.
- Average the multiple summaries and combine within- and between-sample variances. (see Raghunathan 2004)

Multiple Imputation: Create Data Sets

ID	Y(1)	Y(2)	Y(4)	Y(6)	Y(8)
1	93	49	40	46	62
2	79	55	58	65	51
3	70	55	74	NA	NA
4	66	89	72	81	NA
5	73	75	72	62	73
6	61	49	53	NA	NA
7	99	105	98	NA	NA
8	79	59	77	79	134
9	112	141	NA	NA	NA
10	121	NA	NA	NA	NA

Multiple Imputation: Create Data Sets

Data 1					Data 2				
93	49	40	46	62	93	49	40	46	62
79	55	58	65	51	79	55	58	65	51
70	55	74	56	45	70	55	74	65	81
66	89	72	81	81	66	89	72	81	75
73	75	72	62	73	73	75	72	62	73
61	49	53	67	68	61	49	53	24	17
99	105	98	76	91	99	105	98	111	132
79	59	77	79	134	79	59	77	79	134
112	141	118	114	128	112	141	162	141	136
121	116	126	125	130	121	147	125	140	105

Multiple Imputation: Analyze Data Sets

- Linear regression of $Y(8)$ on baseline and treatment group

Imputation	(Intercept)	baseline	Risp (6mg)	Placebo
1	19.00	0.738	-13.32	9.31
2	18.89	0.757	-12.96	9.42
3	15.95	0.790	-12.02	10.55
4	13.43	0.793	-10.94	6.93
5	16.26	0.778	-12.15	10.40

Multiple Imputation: Combine Results

- For each of M imputed data sets you get an estimate of the regression coefficient and the variance of the estimate.
- Average the M estimates.
- Average the M variance estimates and combine with the observed variance of the regression estimates across the imputed data sets to get a final variance (standard error) estimate.

Analysis of PANSS Data

	Final Outcome (8wks)		
	est.	(s.e.)	p-value
Haloperidol is reference group			
Risperidone (6mg)	-7.901	(4.148)	0.0580
Placebo	3.962	(4.844)	0.4142

	Multiple Imputation		
	est.	(s.e.)	p-value
Haloperidol is reference group			
Risperidone (6mg)	-12.278	(3.727)	0.0010
Placebo	9.323	(3.756)	0.0131

Pattern Mixture Models

- Little (1993, 1994, 1995) considers the distribution of \mathbf{Y}_i conditional on the drop-out.
- Define: $D_i = \text{the drop-out time} = \min_k (R_{ik} = 0)$.
- Then we can factor the *complete data* likelihood:

$$\begin{aligned} f(\mathbf{Y}_i, \mathbf{R}_i) &= f(\mathbf{Y}_i, D_i) \\ &= f(\mathbf{Y}_i | D_i) f(D_i) \end{aligned}$$

- This leads to the *observed data* likelihood:

$$\begin{aligned} f(\mathbf{Y}_i^O, \mathbf{R}_i) &= \int_{\mathbf{Y}_i^M} f(\mathbf{Y}_i, D_i) \\ &= \int_{\mathbf{Y}_i^M} f(\mathbf{Y}_i^M | \mathbf{Y}_i^O, D_i) f(\mathbf{Y}_i^O | D_i) f(D_i) d\mathbf{Y} \\ &= f(\mathbf{Y}_i^O | D_i) f(D_i) \end{aligned}$$

- This formulation makes it clear that there's no information in the data about $f(\mathbf{Y}_i^M | \mathbf{Y}_i^O, D_i)$.
- Identifying restrictions.

“Mixing” Pattern Specific Models

Example:

- Suppose that we adopt a linear model specific to the drop-out time:

$$E[Y_{ij} | X_{ij}, D_i = d] = \beta_0^{(d)} + \beta_1^{(d)} X_{ij}$$

- And we have the distribution of drop-out times, $\pi_d = P(D_i = d)$.

Then we can calculate the unconditional expectation:

$$E[Y_{ij} | X_{ij}] = \sum_d \left(\beta_0^{(d)} + \beta_1^{(d)} X_{ij} \right) \cdot \pi_d$$

- If X_{ij} is time, extrapolation $E[Y_{ik} | D_i = d]$ for $k > d$ is clear.
- Dependence of D_i on \mathbf{X}_i may need to be considered.
- Fitzmaurice and Laird (2000)

GEE Modification – Semiparametric Approaches

- Robins, Rotnitzky & Zhao (1995)

★ If the data are MAR and the selection model can be correctly specified, then weighted estimating equations can be constructed that permit valid semi-parametric estimation.

○ Define:

$$\lambda_{ij}(\boldsymbol{\alpha}) = P(R_{ij} = 1 \mid R_{ij-1} = 1, \mathcal{H}_{ij-1}^Y, \mathbf{X}_i)$$

conditional probability of observing Y_{ij}

$$\pi_{ij}(\boldsymbol{\alpha}) = \prod_{k=1}^j \lambda_{ik}(\boldsymbol{\alpha})$$

probability that $R_{ij} = 1$ given \mathcal{H}_{in}^Y

$$\mathcal{H}_{ij-1} = (Y_{i1}, Y_{i2}, \dots, Y_{ij-1})$$

Semiparametric Approaches

○ Define:

$$\mathbf{W}_i = \begin{bmatrix} \frac{R_{i1}}{\pi_{i1}} & 0 & \dots & 0 \\ 0 & \frac{R_{i2}}{\pi_{i2}} & \dots & 0 \\ \vdots & \ddots & & \vdots \\ 0 & 0 & \dots & \frac{R_{in}}{\pi_{in}} \end{bmatrix}$$

● Use an estimating equation of the form

$$U(\boldsymbol{\beta}, \hat{\boldsymbol{\alpha}}) = \sum_{i=1}^N \mathbf{D}_i^*(\mathbf{X}_i, \boldsymbol{\beta}) \mathbf{W}_i(\hat{\boldsymbol{\alpha}}) [\mathbf{Y}_i - \boldsymbol{\mu}_i(\boldsymbol{\beta})]$$

- This is an unbiased estimating equation because:

$$E \left[\frac{R_{ij}}{\pi_{ij}} (Y_{ij} - \mu_{ij}) \right] = 0$$

- The matrix $\mathbf{D}_i^*(\boldsymbol{\beta})$ can be chosen to obtain efficiency.

★★ See Preisser, Lohman & Rathouz (2002) for study of IPW.

GEE and Conditional Mean Imputation

- Paik (1997) considers use of standard GEE estimators when data are MAR by adding a conditional mean model for imputation of missing data.

- Define:

$$\mathbf{Y}_i^* = \mathbf{R}_i \cdot \mathbf{Y}_i + (\mathbf{1} - \mathbf{R}_i) \tilde{\mathbf{Y}}_i$$

$$\tilde{Y}_{ij} = E[Y_{ij} \mid \mathcal{H}_{ij-1}^Y, \mathbf{X}_i]$$

- Replace missing outcomes using the mean of subjects with the same history, \mathcal{H}_{ij-1}^Y .
- Complete $\tilde{\mathbf{Y}}_i$ obtained via sequential imputation.
- Multiple imputation.

Summary

- Missing data can bias results.
- Study reasons for missingness.
- Likelihood and Multiple Imputation methods may help.
- Results depend on assumptions.
- Conduct sensitivity analyses.