



COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK

STAT 4224/5224

Bayesian Statistics

Dobrin Marchev

Missing Data

- Missing values: aka most common type of data
- In R they are denoted NA
- Why do we care?

Simple example in R:

```
x = c(4, 6, 2)
```

```
mean(x)
```

```
[1] 4
```

```
x = c(4, 6, NA)
```

```
mean(x)
```

```
[1] NA
```

Overview of Methods

Methods that throw away data (older approaches, prior to 1987)

- Complete cases analysis aka CCA (not a great idea)
- Complete variables (bad idea)
- Pairwise deletion (for special purposes)

Methods that don't throw away data (new approaches, 1990s to present)

- Imputation: single and multiple

CCA

X_1	X_2	X_3	X_4
0	2	5	2
0	1	3	NA
0	2	4	3
0	3	5	NA
1	5	NA	NA
1	4	NA	12
1	6	11	NA
1	6	12	16



X_1	X_2	X_3	X_4
0	2	5	2
0	2	4	3
1	6	12	16

Complete cases (aka listwise deletion)

- Removes all observations from the dataset that have any missing values (most software packages do this automatically when you run an analysis, for example, the `lm` function in R)
- At best it is inefficient (yields higher standard errors) because of reduced sample size
- At worst it can cause severe bias
- Reductions in sample size may preclude certain types of analyses, e.g. subgroup analyses
- When does it work?

Missing Data Mechanisms

Rubin (1976) classified missing data into 3 categories:

- Missing Completely at Random (MCAR)
- Missing at Random (MAR)
- Not Missing at Random (NMAR), also called Missing Not at Random (MNAR)
- Aka the most confusing statistical terms ever invented

The process that governs the probability of a data point being missing is called *missing data mechanism* or response mechanism.

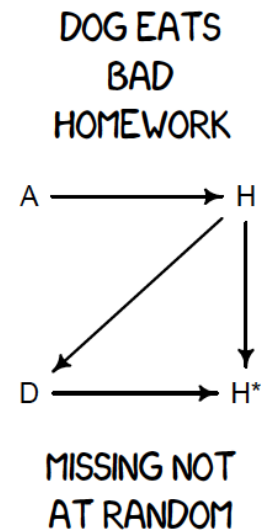
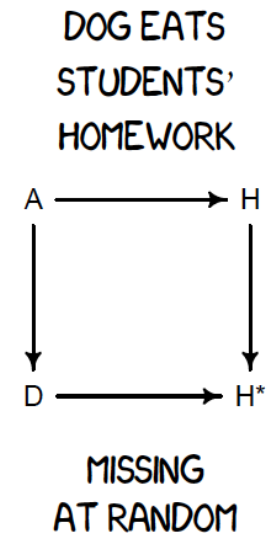
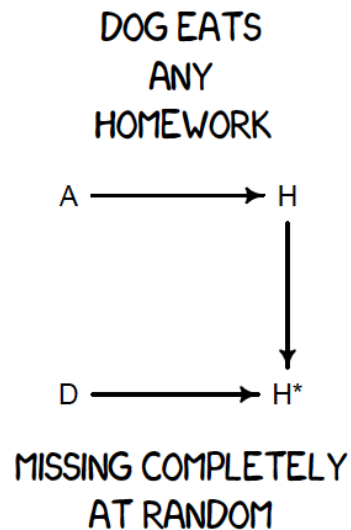
MDM: Example

Consider a cancer study in which one of the variables is a quality-of-life score, which has missing entries.

- If missing entries are due to a random computer glitch or human error, unrelated to the patients at all, then we have MCAR
- If missing quality of life is a function of age and education, or treatment and health status, which we observe directly, then we have MAR
- If missing is direct function of quality of life, then we have NMAR

- H: Homework
- H*: Homework with missing values
- A: Attribute of student (say, study time)
- D: Dog (missingness mechanism)

Three Types of Missingness



Notation

- Sample data matrix is usually denoted \mathbf{X} , that is, the $n \times p$ matrix containing the data values on p variables for all n units in the sample.
- We define the *response indicator* \mathbf{R} as an $n \times p$ 0–1 matrix (see next slide for details).
- Specific elements of \mathbf{X} and \mathbf{R} are denoted by y_{ij} and r_{ij} , respectively.
- We are restricted to the case where \mathbf{R} is completely known, i.e., we know where the missing data are. This covers many applications of practical interest, but not all. For example, some questionnaires present a list of diseases and ask the respondent to place a “tick” at each disease that applies. If there is a “yes” we know that the field is not missing. However, if the field is not ticked, it could be because the person didn’t have the disease (a genuine “no”) or because the respondent skipped the question (a missing value).

Notation

Let R be the **matrix** of variables R_1, \dots, R_p , corresponding to the variables in our dataset, X_1, \dots, X_p , that indicate whether a given value of the corresponding X variable is observed ($= 1$) or missing ($= 0$)

X_1	X_2	X_3	X_4
0	2	5	2
0	1	3	?
0	2	4	3
0	3	5	?
1	5	?	?
1	4	?	12
1	6	11	?
1	6	12	16

R_1	R_2	R_3	R_4
1	1	1	1
1	1	1	0
1	1	1	1
1	1	1	0
1	1	0	0
1	1	0	1
1	1	1	0
1	1	1	1

Missing Completely at Random

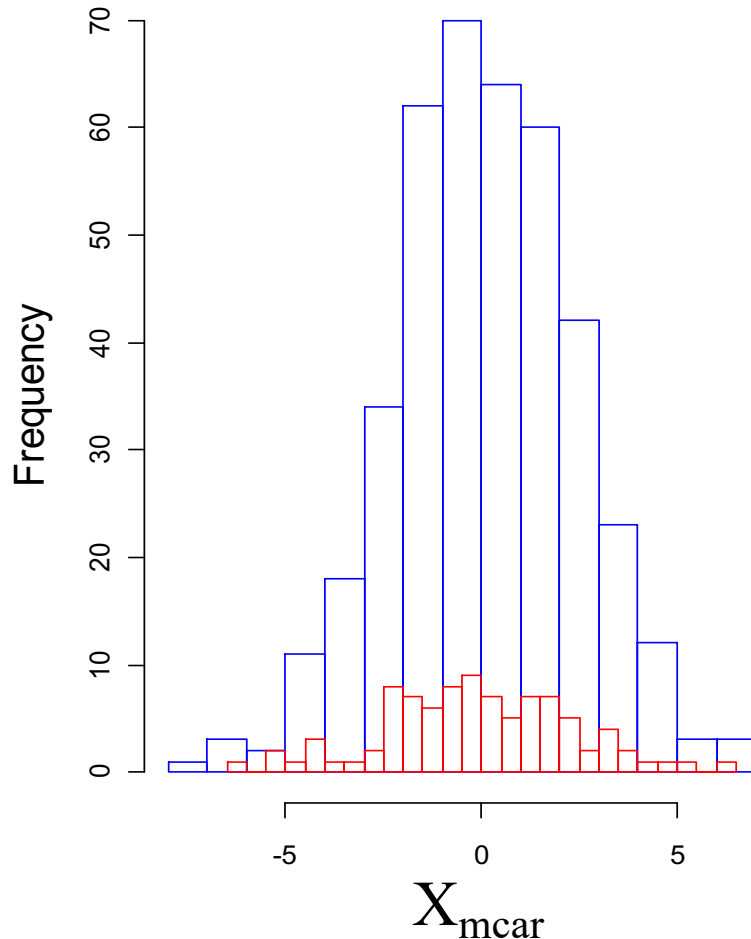
- More formally, let \mathbf{R} be the missing pattern and \mathbf{X} the data
- In MCAR \mathbf{R} and \mathbf{X} are assumed independent random vectors ($\mathbf{R} \perp \mathbf{X}$)
- $P(R_1, R_2, \dots, R_p \mid X_1, X_2, \dots, X_p) = P(R_1, R_2, \dots, R_p)$
- $P(\mathbf{R} = \mathbf{0} \mid \mathbf{X}_{\text{obs}}, \mathbf{X}_{\text{mis}}, \varphi) = P(\mathbf{R} = \mathbf{0} \mid \varphi)$
- This means that whether any given value is missing is *completely random* (doesn't depend on any other variable).
- This is generally not a plausible assumption. Usually, certain types of people/patients are much more likely than others to have missing data.
- When could this happen in reality?
 - A bunch of records are lost
 - Missing by design

MCAR data (if we could observe it!):

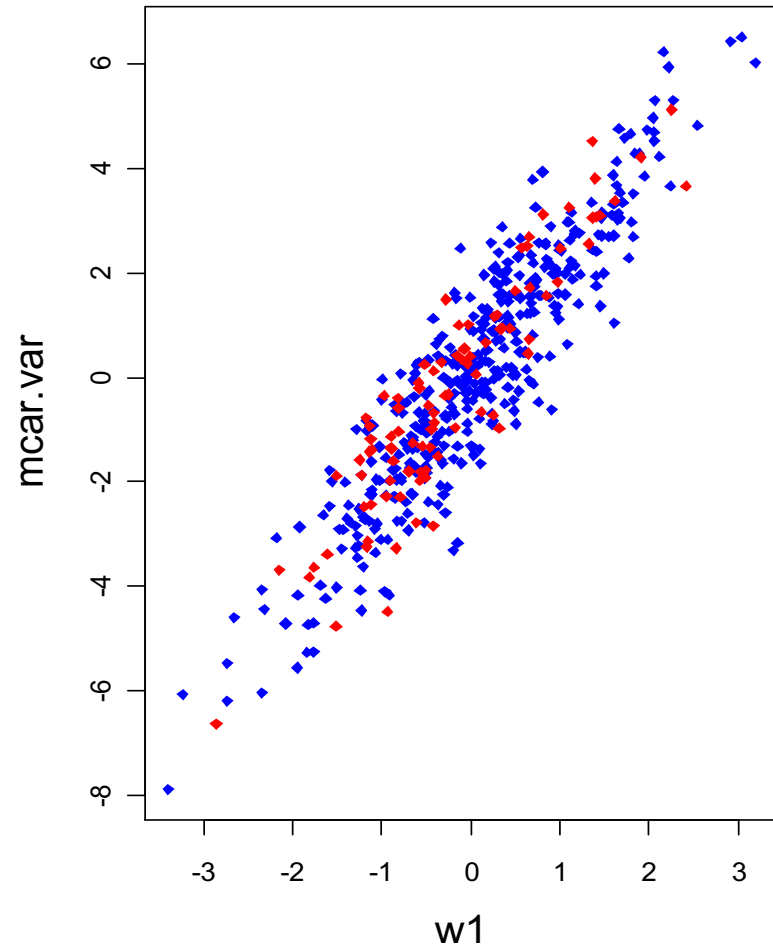
blue = observed

red = missing

observed and missing data



observed and missing data



Estimation with MCAR

- Suppose we are interested in estimating the mean of the population $\mu = E(X)$, but instead of iid sample X_1, \dots, X_n , we observe X_i if $R_i = 1$, and X_i is missing when $R_i = 0$.
- The complete cases estimator is:

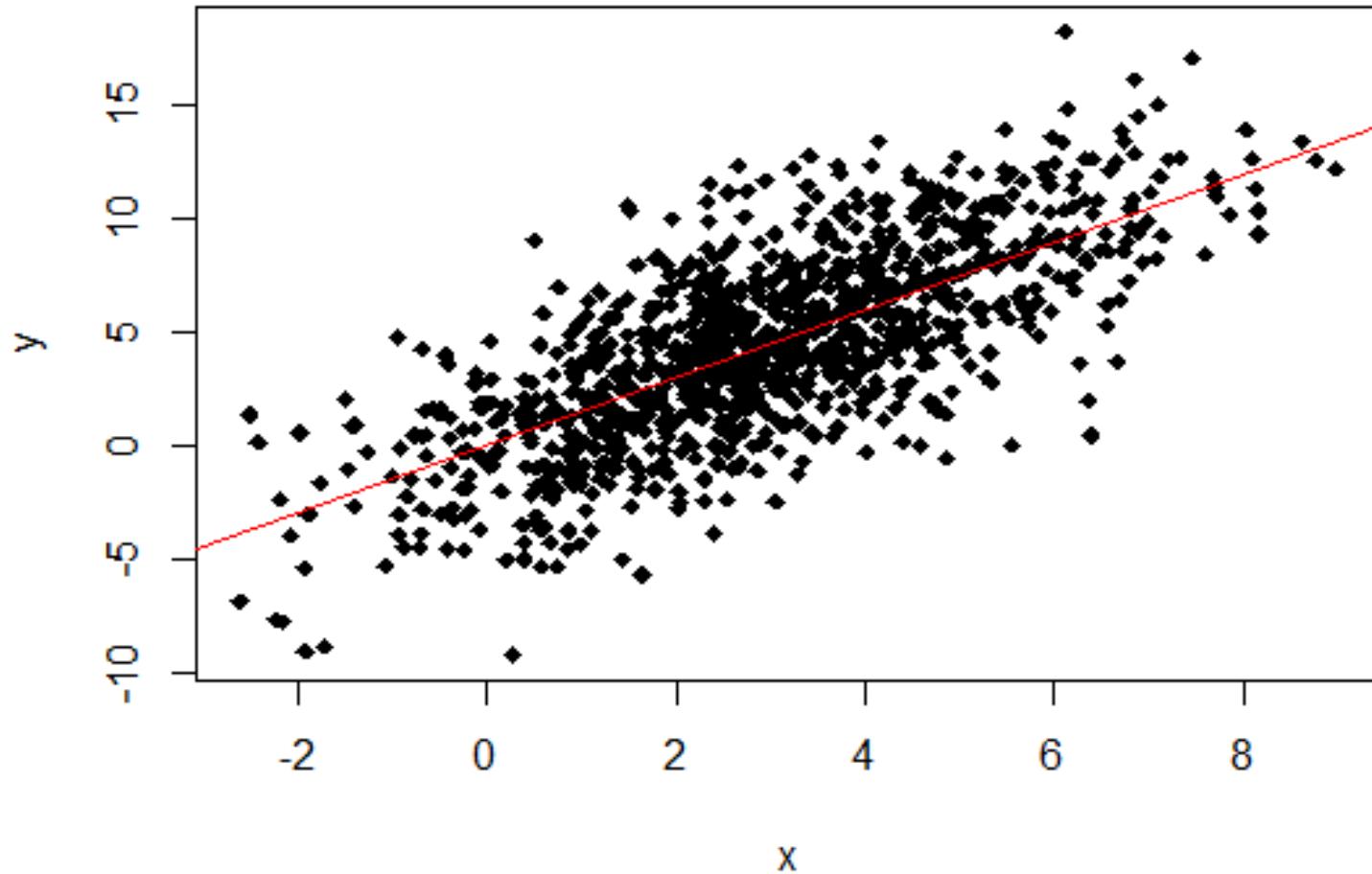
$$\hat{\mu}^C = \frac{\sum_{i=1}^n R_i X_i}{\sum_{i=1}^n R_i}$$

- Assuming MCAR, that is, $P(R = 1 | X) = P(R = 1)$, it can be shown that

$$\hat{\mu}^C = \frac{\frac{1}{n} \sum_{i=1}^n R_i X_i}{\frac{1}{n} \sum_{i=1}^n R_i} \xrightarrow{P} \frac{E(RX)}{E(R)} = \frac{E(R)E(X)}{E(R)} = E(Y) = \mu$$

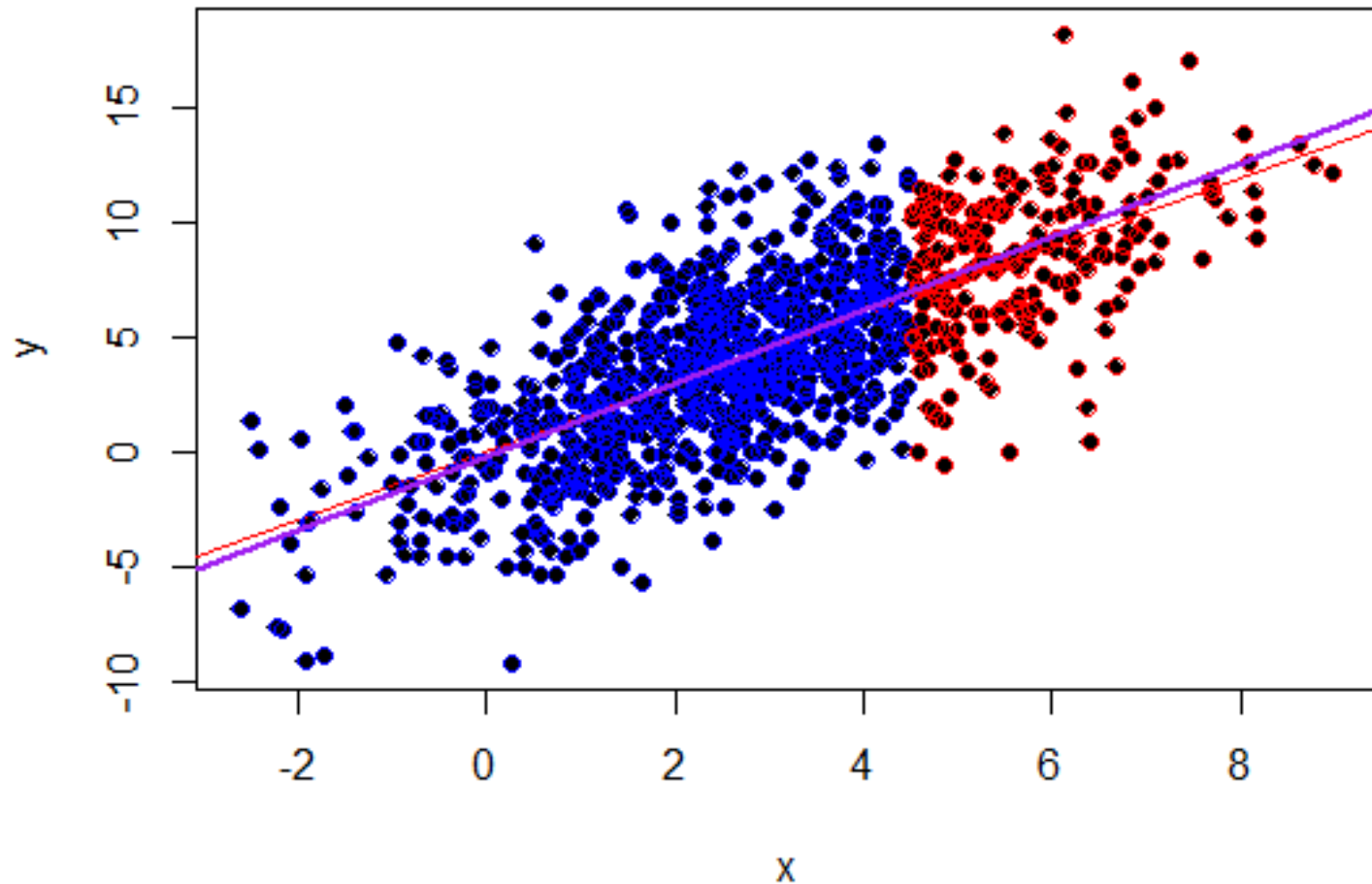
Illustration of CCA and

.



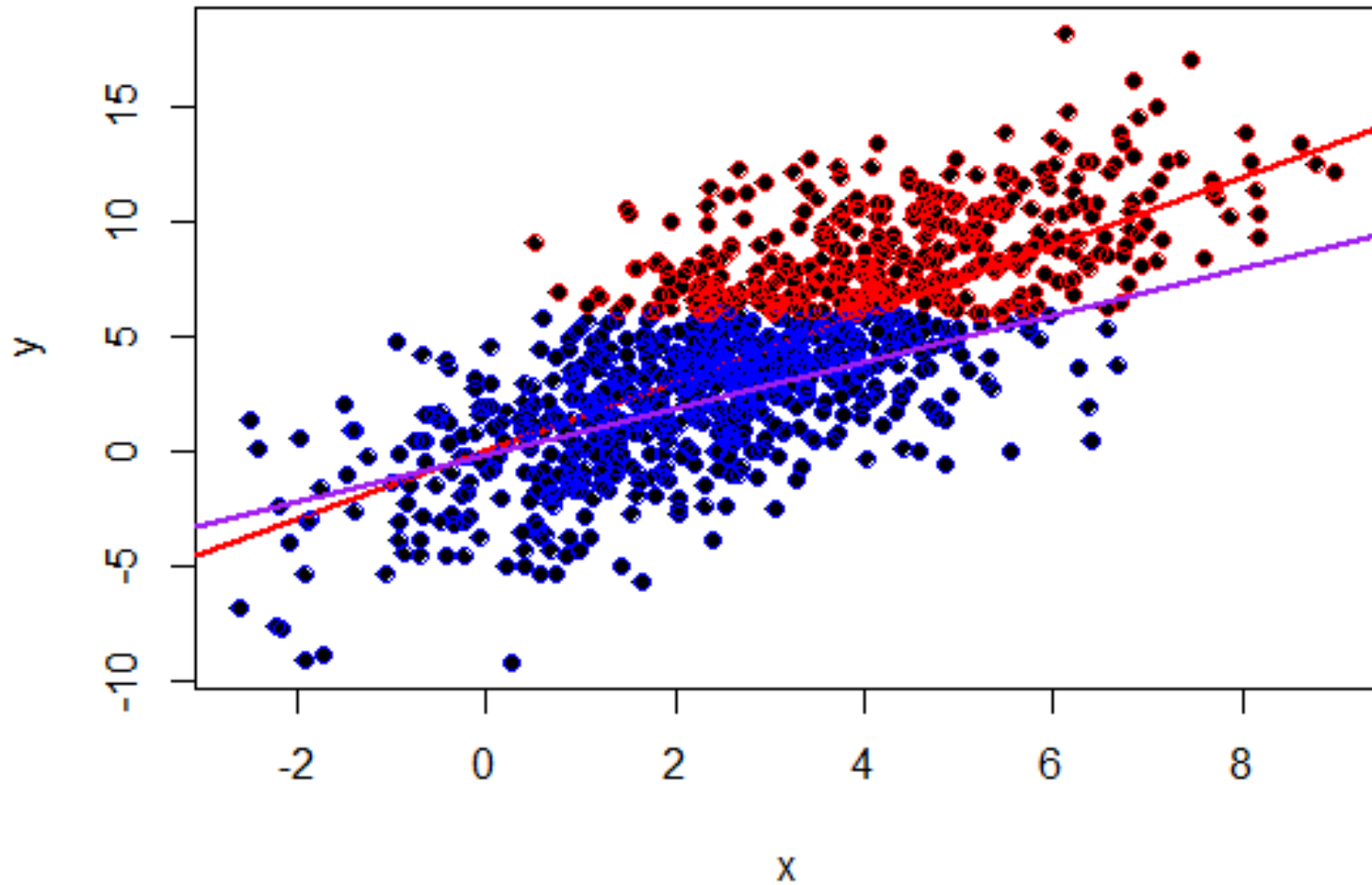
CCA and regression: just fitting to observed (blue) points

Missingness due to x



CCA and regression: just fitting to observed (blue) points

Missingness due to y



Pairwise Deletion

- Procedure that focuses on the covariance matrix.
- Each element of that matrix is estimated from all data available for that element.
- Because different variances and covariances are based on different subsamples of respondents, parameter estimates may be biased unless missingness is MCAR.
- (More technical): because the different parameters are estimated with different subsamples, it often happens that the matrix is not positive definite!

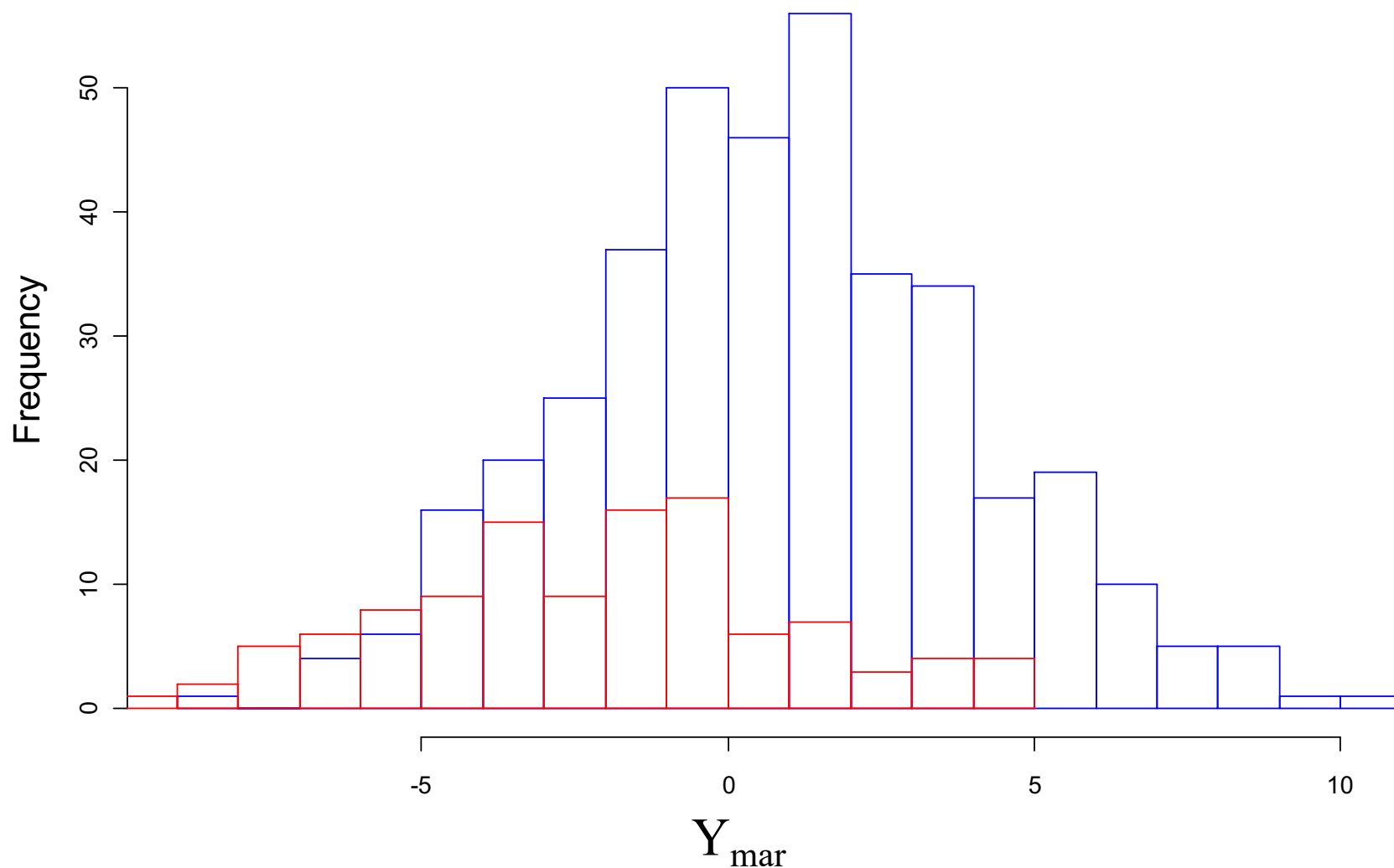
Missing at Random (MAR)

- The propensity for a data point to be missing is not related to the missing data, but it is related to some of the observed data.
- $P(R_1, \dots, R_p \mid X_1, \dots, X_p) = P(R_1, \dots, R_p \mid X_1^{\text{obs}}, \dots, X_p^{\text{obs}})$
- $P(\mathbf{R} \mid \mathbf{X}) = P(\mathbf{R} \mid \mathbf{X}^{\text{obs}})$
- Here missingness depends on observed values of the variables.
- A simple version of this is $P(R_1 \mid X_1, \mathbf{W}) = P(R_1 \mid \mathbf{W})$, where \mathbf{W} is a subset of fully observed variables in the matrix \mathbf{X}
- Classic example:

With a long, self-administered survey, for which there is a limited amount of time for completion, fast readers will complete the survey, but slow readers will leave some questions blank at the end. However, reading speed is something that can be measured early in the questionnaire where virtually all of the respondents will provide data.

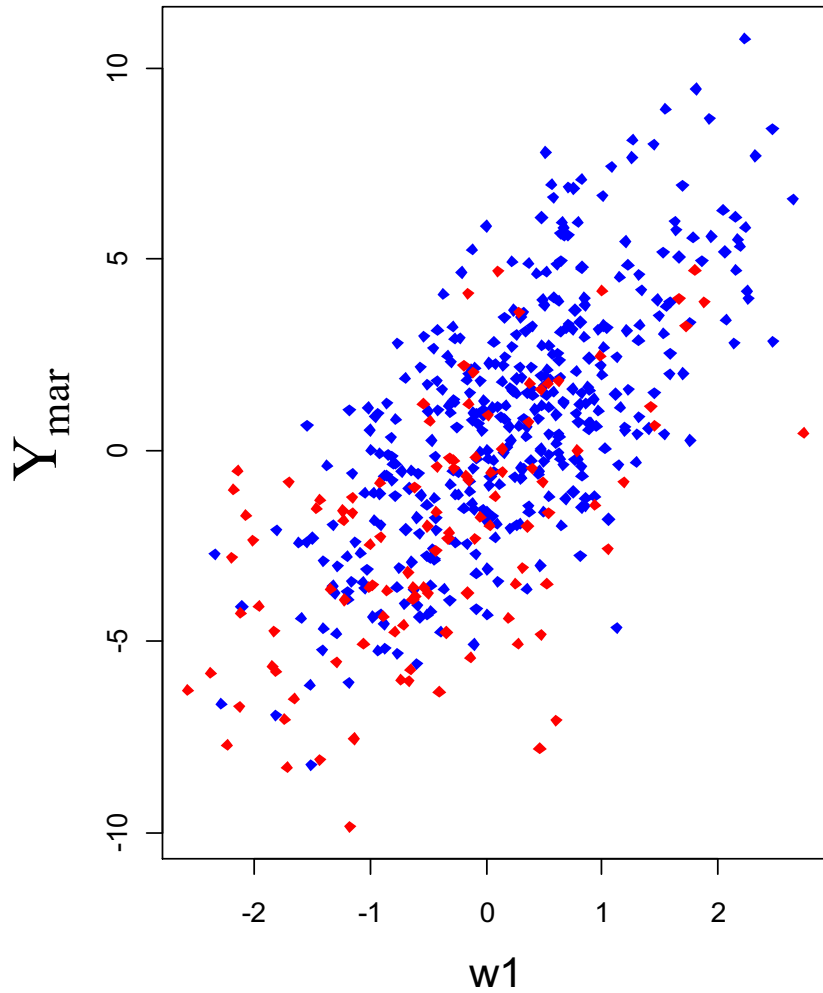
MAR data (if we could see everything!):

observed and missing data

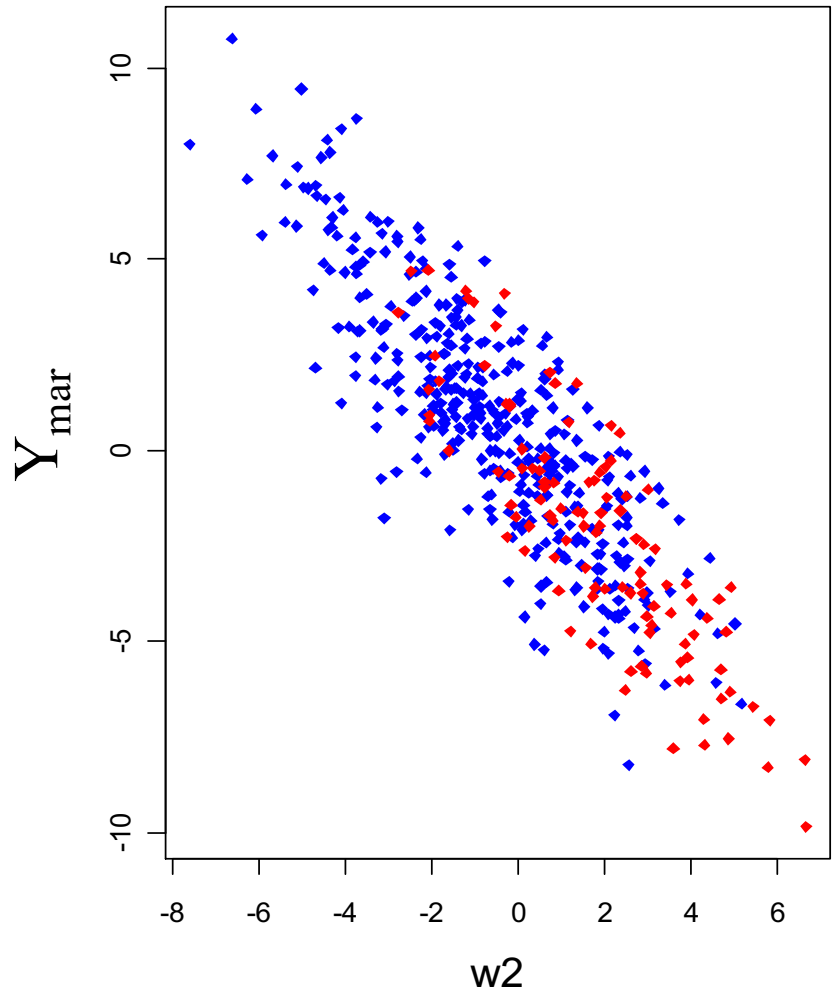


MAR data:

observed and missing data

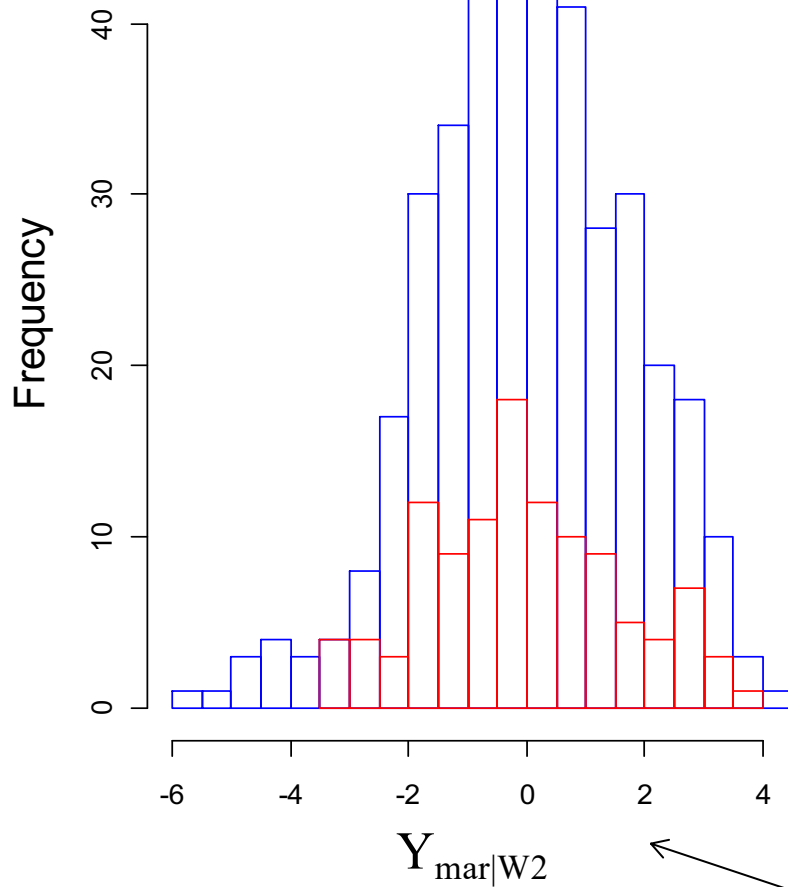


observed and missing data

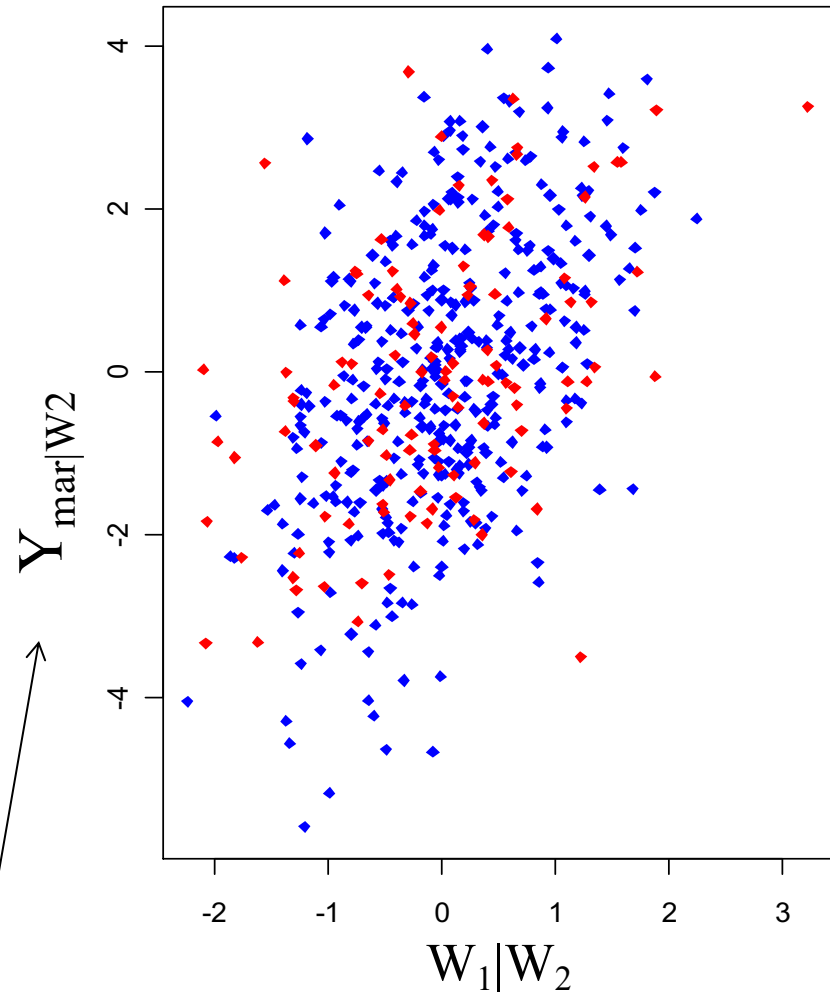


MAR data: Y_{mar} *conditional* on W_2

observed and missing data



observed and missing data



Where $Y_{\text{mar}}|W_2$ is the residuals from a regression of Y_{mar} on W_2

Not Missing at Random (NMAR)

- In this case missingness depends on the values of the items that are missing!

- Thus

$$P(R_1, \dots, R_p \mid Y_1, \dots, Y_p) \neq P(R_1, \dots, R_p \mid Y_1^{\text{obs}}, \dots, Y_p^{\text{obs}})$$

$$P(\mathbf{R} \mid \mathbf{Y}) \neq P(\mathbf{R} \mid \mathbf{Y}^{\text{obs}})$$

- One way of formalizing this is: $P(\mathbf{R} \mid \mathbf{Y}) = P(\mathbf{R} \mid \mathbf{Y}^{\text{obs}}, \mathbf{Z})$
- That is, NMAR missingness occurs when missingness on \mathbf{Y} (i.e., \mathbf{R}) is caused by \mathbf{Y} itself, by some variant of \mathbf{Y} , or by some other variable that is related to \mathbf{Y} , but which has not been measured.
- Example: very wealthy people are less likely to report their income *and this wealth is not predicted by the other variables in the data*

Recall: Multivariate Normal Model

Assume that we have multivariate observations

$$\mathbf{X}_1, \dots, \mathbf{X}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma} \sim N_p(\boldsymbol{\theta}, \boldsymbol{\Sigma})$$

If there are no missing values, then the likelihood is

$$\begin{aligned} f(\mathbf{x}_1, \dots, \mathbf{x}_n \mid \boldsymbol{\theta}, \boldsymbol{\Sigma}) \\ &= \prod_{i=1}^n (2\pi)^{-\frac{p}{2}} (\det \boldsymbol{\Sigma})^{-1/2} e^{-1/2 (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})} \\ &= (2\pi)^{-\frac{np}{2}} (\det \boldsymbol{\Sigma})^{-\frac{n}{2}} e^{-1/2 \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})} \end{aligned}$$

Q: How do we compute $f(\mathbf{x}_i \mid \boldsymbol{\theta}, \boldsymbol{\Sigma}) =$

$(2\pi)^{-\frac{p}{2}} (\det \boldsymbol{\Sigma})^{-1/2} e^{-1/2 (\mathbf{x}_i - \boldsymbol{\theta})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\theta})}$ when the \mathbf{x}_i is missing?

A: Impute a value.

Note: The worst idea is to impute the average!

Likelihood with missing data

Assume MAR. Then the likelihood of observed data for subject i is

$$\begin{aligned} f(\mathbf{r}_i, \{x_{ij} : r_{ij} = 1\} | \boldsymbol{\theta}, \boldsymbol{\Sigma}) &= p(\mathbf{r}_i) \times f(\{x_{ij} : r_{ij} = 1\} | \boldsymbol{\theta}, \boldsymbol{\Sigma}) \\ &= p(\mathbf{r}_i) \times \int \left[f(\mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\Sigma}) \prod_{J:r_{ij}=0} dx_{ij} \right] \end{aligned}$$

That is, we multiply $p(\mathbf{r}_i)$ with the marginal distribution of the observed variables, after integrating out the missing variables.

For example, if $\mathbf{x}_i = (x_{i1}, NA, NA, x_{i4})$, then $\mathbf{r}_i = (1, 0, 0, 1)$ and

$$f(\mathbf{r}_i, x_{i1}, x_{i4} | \boldsymbol{\theta}, \boldsymbol{\Sigma}) = p(\mathbf{r}_i) \times \iint f(\mathbf{x}_i | \boldsymbol{\theta}, \boldsymbol{\Sigma}) dx_{i2} dx_{i3}$$

So, the correct thing to do when data are missing at random is to integrate over the missing data to obtain the marginal probability of the observed data.

Algorithm

- Let \mathbf{X} be the $n \times p$ matrix of all data, both observed and missing.
- Let \mathbf{R} be the $n \times p$ missing pattern matrix as defined before.
- The matrix \mathbf{X} consists of two parts:
- $\mathbf{X}_{\text{obs}} = \{x_{ij} : r_{ij} = 1\}$ is the observed data.
- $\mathbf{X}_{\text{mis}} = \{x_{ij} : r_{ij} = 0\}$ is the unobserved or missing data.
- Note that \mathbf{X}_{mis} has to be treated as unknown parameter!
- Goal: Obtain samples from the posterior distribution
$$f(\boldsymbol{\theta}, \boldsymbol{\Sigma}, \mathbf{X}_{\text{mis}} | \mathbf{X}_{\text{obs}})$$
- Solution: Run a three-stage Gibbs sample that iterates between
$$f(\boldsymbol{\theta} | \mathbf{X}_{\text{obs}}, \boldsymbol{\Sigma}, \mathbf{X}_{\text{mis}})$$
$$f(\boldsymbol{\Sigma} | \mathbf{X}_{\text{obs}}, \boldsymbol{\theta}, \mathbf{X}_{\text{mis}})$$
$$f(\mathbf{X}_{\text{mis}} | \mathbf{X}_{\text{obs}}, \boldsymbol{\theta}, \boldsymbol{\Sigma})$$