# MISSING DATA

BST228

# THE MISSING DATA PROBLEM

We want to learn about parameters $\boldsymbol{\theta}$ from observed data $\mathbf{y}_{\mathrm{obs}}$, acknowledging missing data $\mathbf{y}_{\mathrm{mis}}$. The posterior is

$$p\left(\boldsymbol{\theta} \mid \mathbf{y}_{\mathrm{obs}}\right) \propto p\left(\mathbf{y}_{\mathrm{obs}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right).$$

**How do we account for $\mathbf{y}_{\mathrm{mis}}$?**

- Sometimes we can't or choose to observe all the data, e.g., due to respondents not answering survey questions, budget constraints, or individuals dropping out of clinical studies.
- Handling missing data is all about still making valid inferences in light of this challenge.

# MISSING DATA AS PARAMETERS

We treat missing data $\mathbf{y}_{\mathrm{mis}}$ like any other parameter

$$p\left(\boldsymbol{\theta}, \mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}\right) \propto p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right).$$

# MISSING DATA COMPONENTS

1. Complete data likelihood $p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right)$.
2. Observation indicator $\mathbf{I}$ for each data point; 1 if observed, 0 if not.
3. Missingness model $p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}}, \mathbf{x}, \boldsymbol{\phi}\right)$ with parameters $\boldsymbol{\phi}$.
4. Joint posterior $p\left(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right)$.

Speaker notes

- The complete data likelihood $p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right)$ is the likelihood we could evaluate if we had both the observed data $\mathbf{y}_{\mathrm{obs}}$ and the missing data $\mathbf{y}_{\mathrm{mis}}$.
- We treat the observation indicator $\mathbf{I}$ as fully observed—we generally know if we have the data or not.
- The missingness model captures how and why the data might be present or absent.
- If we can sample from the joint posterior, we can infer the parameters $\boldsymbol{\theta}$ we are interested in and ignore the imputed missing data $\mathbf{y}_{\mathrm{mis}}$ and parameters of the $\boldsymbol{\phi}$ of the missingness model.

# MISSING DATA MECHANISMS

- Missing Completely at Random (MCAR): Data are missing independent of $\mathbf{y}_{\mathrm{obs}}$ and $\mathbf{y}_{\mathrm{mis}}$.
- Missing at Random (MAR): Conditional on $\mathbf{y}_{\mathrm{obs}}$, data are missing independent of $\mathbf{y}_{\mathrm{mis}}$.
- Missing not at Random (MNAR): Missingness depends on $\mathbf{y}_{\mathrm{mis}}$.

Speaker notes

- We consider the same dataset (1,000 test scores in the first year and 800 test scores in the second year) but with different scenarios for why the data are missing.

# EXAMPLE: TEST SCORES

Test scores for 1,000 students will be measured for two years, but 200 values are missing in the second year.

**A**: Students who scored lower in the first year are less likely to accept the invitation to take the test in the second year.

# EXAMPLE: TEST SCORES

Test scores for 1,000 students will be measured for two years, but 200 values are missing in the second year.

**B**: Budget constraints in the second year require reducing the sample size. A random subset comprising 800 students are selected to take the test.

# EXAMPLE: TEST SCORES

Test scores for 1,000 students will be measured for two years, but 200 values are missing in the second year.

**C**: Students who feel they didn't do well on the test do not hand in their work.

# RESULTS FROM POLL FOR MISSINGNESS MECHANISMS

| Scenario | MCAR | MAR | MNAR |
| --- | --- | --- | --- |
| A: not returning for second test | 8% | **62%** | 31% |
| B: budget constraints | **85%** | 8% | 8% |
| C: not handing in poor results | 0% | 15% | **85%** |

Speaker notes

- Our online poll correctly identifies the missing data mechanisms.
- Students not returning for the second test is data MAR because the missingness depends on the observed test scores from the previous year.
- Budget constraints lead to data MCAR because the budget has nothing to do with the test scores (unless the budget affects teaching).
- Not handing in poor results leads to data MNAR because the missingness depends on the data we are trying to collect.

# THE DATA ARE …

- **A**: missing at random, $p\left(\mathbf{I} \mid \mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}}, \phi\right) = p\left(\mathbf{I} \mid \mathbf{y}_{\text{obs}}, \phi\right)$
- **B**: missing completely at random, $p\left(\mathbf{I} \mid \mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}}, \phi\right) = p\left(\mathbf{I} \mid \phi\right)$
- **C**: missing not at random, $p\left(\mathbf{I} \mid \mathbf{y}_{\text{obs}}, \mathbf{y}_{\text{mis}}, \phi\right) = p\left(\mathbf{I} \mid \mathbf{y}_{\text{mis}}, \phi\right)$

# WHY DOES THIS MATTER?

We use the following random effects model for test scores:

$$\mu_t \sim \text{Normal}\,(0, 1) \quad \text{for } t \in \{1, 2\}$$

$$a_i \sim \text{Normal}\,(0, \kappa^2) \quad \text{for } i \in \{1, \ldots, 1000\}$$

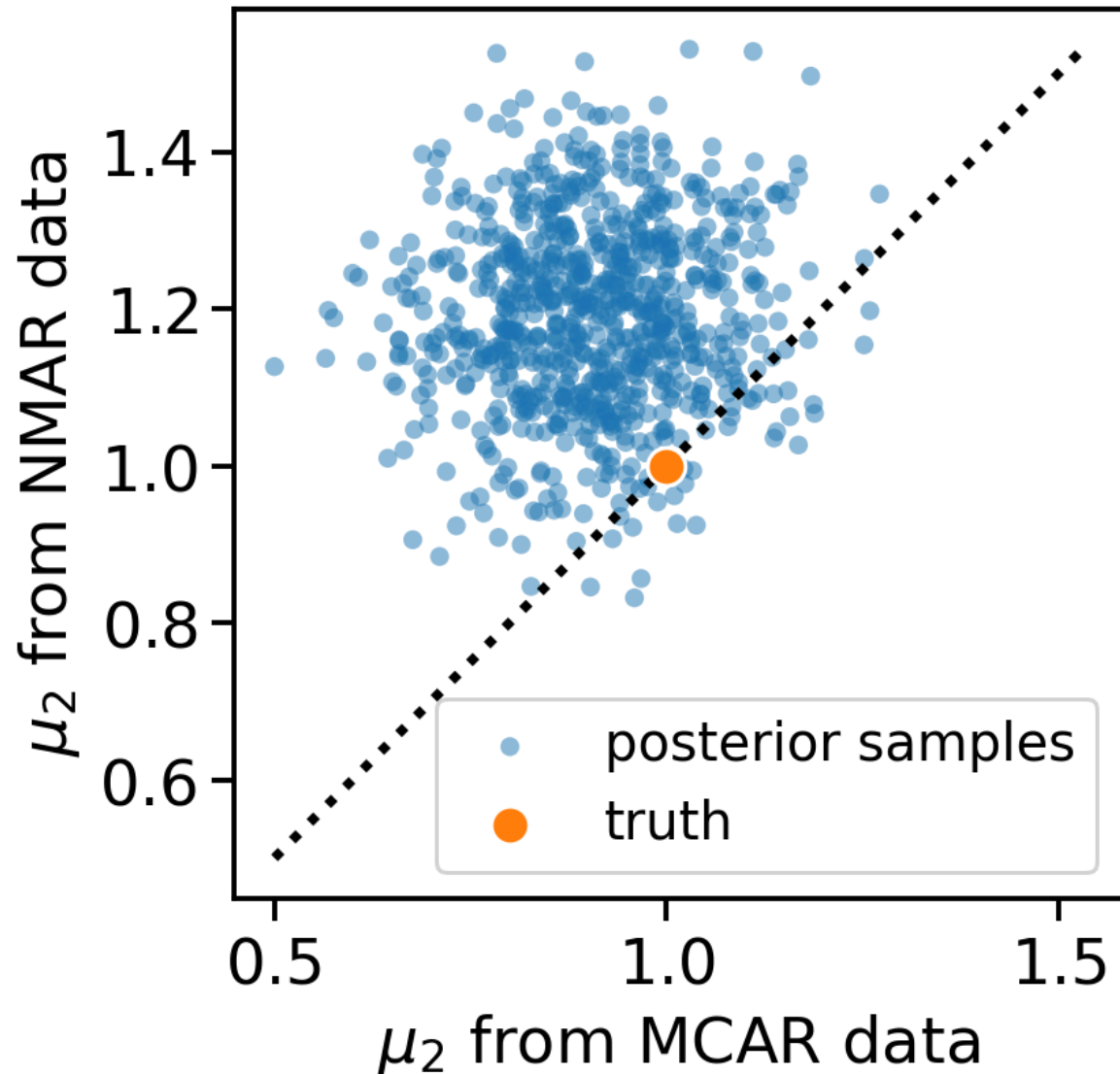$$y_{it} \sim \text{Normal}\,(a_i + \mu_t, \sigma^2),$$

where $\mu$ is the mean score for each year and $a_i$ is the skill of each student.

Speaker notes

- Test scores $y_{it}$ depend both on the overall score $\mu_t$ for the test in year $t$ (a measure of how easy the test is) and the individual student effect $a_i$ for student $i$ (a measure of how skilled the student is at answering test questions).
- The scale parameter $\kappa$ captures the variance in student abilities. If $\kappa$ is small, all students get similar scores. If $\kappa$ is large, there is substantial variation in ability.

**C**: Students who feel they didn't do well on the test do not hand in their work.
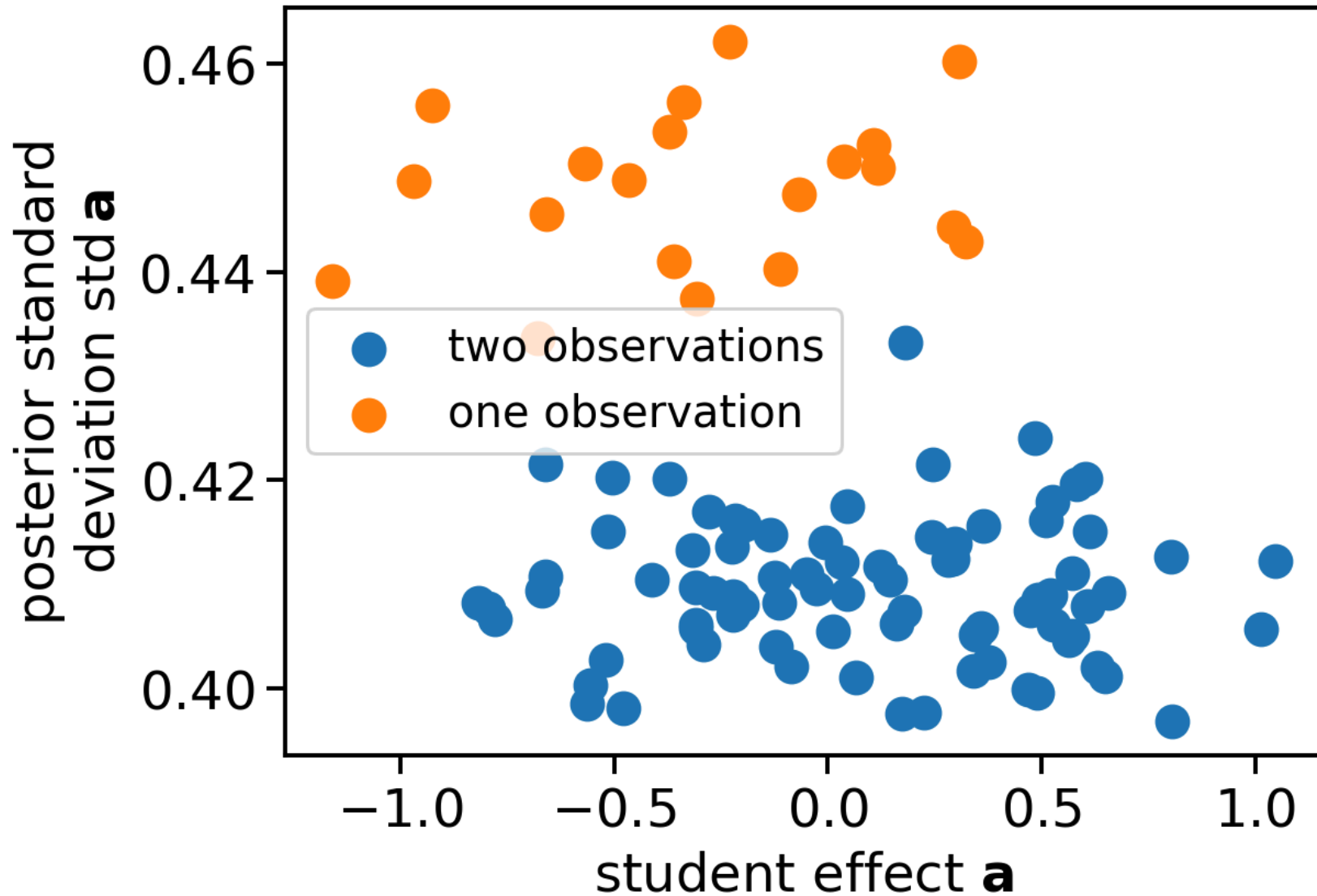
- Blue scatter are posterior samples for a naive model that does not account for the missingness. Samples along the horizontal axis come from scenario B (MCAR due to budget constraints) and samples along the vertical axis come from scenario C (NMAR due to students not handing in results).
- Mean score $\mu_2$ for the second year is overestimated because low scores are missing from the dataset.

**A**: Students who scored lower in the first year are less likely to accept the invitation to take the test in the second year.

- The figure shows a scatter of posterior standard deviation for the inferred student effects $a$ against the true student effect (which we know because the data are simulated).
- We make two observations: (a) are larger for students who did not complete the second test. (b) students who did not complete the second test have lower ability.
- Unlike in the NMAR setting, results for both $\mu$ and $\mathbf{a}$ are unbiased. Intuitively, that's because row means of $\mathbf{y}$ (where entries are available) are unbiased estimators of student ability and columns means of $\mathbf{y}$ (where entries are available) are unbiased estimators of overall exam score.

# JOINT VS MARGINAL POSTERIOR

We can sample …

- from the joint posterior

$$p\left(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right) \propto p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}}, \boldsymbol{\phi}\right)$$
$$\times p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}, \boldsymbol{\phi}\right)$$

- or from the marginal posterior

$$p\left(\boldsymbol{\theta} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right) = \int d\boldsymbol{\phi} \, d\mathbf{y}_{\mathrm{mis}} \, p\left(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right).$$

# IGNORABILITY FOR DATA MAR

Suppose data are MAR such that $p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}}, \boldsymbol{\phi}\right) = p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\phi}\right)$ and the parameters of our model $\boldsymbol{\theta}$ and of the missingness model $\boldsymbol{\phi}$ are independent a priori, i.e., $p\left(\boldsymbol{\theta}, \boldsymbol{\phi}\right) = p\left(\boldsymbol{\theta}\right) p\left(\boldsymbol{\phi}\right)$. Then, using the expression for the joint posterior from the previous slide,

$$p\left(\boldsymbol{\theta} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right) = \int d\boldsymbol{\phi}\, d\mathbf{y}_{\mathrm{mis}}\, p\left(\boldsymbol{\theta}, \boldsymbol{\phi}, \mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right)$$

$$\propto \int d\boldsymbol{\phi}\, d\mathbf{y}_{\mathrm{mis}}\, p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\phi}\right) p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right) p\left(\boldsymbol{\phi}\right).$$

The integral with respect to $\mathbf{y}_{\mathrm{mis}}$ and $\boldsymbol{\phi}$ is separable such that

$$p\left(\boldsymbol{\theta} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{I}\right) \propto \left[\int d\boldsymbol{\phi}\, p\left(\mathbf{I} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\phi}\right) p\left(\boldsymbol{\phi}\right)\right] \left[\int d\mathbf{y}_{\mathrm{mis}}\, p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right)\right]$$

$$\propto \int d\mathbf{y}_{\mathrm{mis}}\, p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right)$$

$$\propto p\left(\mathbf{y}_{\mathrm{obs}} \mid \boldsymbol{\theta}\right) p\left(\boldsymbol{\theta}\right).$$

The second line follows because the first integral does not depend on $\boldsymbol{\theta}$ and can be absorbed by the proportionality. The third line follows by the law of total probability.

The missing data model does not appear in the posterior, and we say the missing data mechanism is *ignorable*. This means we do not need to model $p\left(\mathbf{I} \mid \ldots\right)$, but it does *not* mean that we can ignore that data are missing.

# EVALUATING THE OBSERVED-DATA LIKELIHOOD

Assuming an ignorable missingness mechanism, we seek to evaluate the observed data likelihood

$$p\left(\mathbf{y}_{\mathrm{obs}} \mid \boldsymbol{\theta}\right) = \int d\mathbf{y}_{\mathrm{mis}}\, p\left(\mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}} \mid \boldsymbol{\theta}\right).$$

# DATA AUGMENTATION IN THEORY

If the integral with respect to $\mathbf{y}_{\mathrm{mis}}$ is not tractable, we treat $\mathbf{y}_{\mathrm{mis}}$ as a latent variable and iteratively …

1. sample parameters $\boldsymbol{\theta} \mid \mathbf{y}_{\mathrm{obs}}, \mathbf{y}_{\mathrm{mis}}$ from the posterior given complete data, including imputations $\mathbf{y}_{\mathrm{mis}}$.

2. impute missing data $\mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\theta}$ given observed data and model parameters.

3. repeat, starting at 1.

Speaker notes

- For non-tractable marginalization, we can sample from the posterior predictive distribution $\mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\theta}$ to obtain an "imputed" dataset, i.e., a dataset that is consistent with our model.
- We then sample parameters $\boldsymbol{\theta}$ treating $\mathbf{y}_{\mathrm{mis}}$ as if it was observed.
- What algorithm is this? It is a Gibbs sampler because we iteratively sample from the conditionals.

# DATA AUGMENTATION IN PRACTICE

- Pick initial parameters $\boldsymbol{\theta}$, e.g., based on a complete case analysis.
- Pick initial imputations $\mathbf{y}_{\mathrm{mis}}$, e.g., based on mean imputation.
- Run the sampler.
- Check for convergence using MCMC diagnostics.
- Use $\boldsymbol{\theta}$ samples for analysis and posterior predictions.
- Use $\mathbf{y}_{\mathrm{mis}}$ samples to investigate missing data.

Speaker notes

- In practice, we need to initialize the algorithm with sensible starting values. For parameters $\boldsymbol{\theta}$, this can be achieved by using a complete case analysis (discarding all records that do not have complete data). For missing data $\mathbf{y}_{\mathrm{mis}}$, we can use a simple imputation method like replacing all missing values by the mean of observed values.
- After running the sampler, we need to check for convergence and discard burn-in samples because the sampler needs to converge to the target distribution from the initial guesses.

# MISSING DATA ALWAYS REQUIRES A MODEL

To sample from $\mathbf{y}_{\mathrm{mis}} \mid \mathbf{y}_{\mathrm{obs}}, \boldsymbol{\theta}$, we must have a model for all missing parts, e.g., missing covariates in a regression model.

# EXAMPLE: DIABETES

A diabetes dataset comprises 442 observations of disease progression $\mathbf{y}$ with a design matrix $\mathbf{X}$ comprising ten features for each observation. We assume 10% of predictors are MCAR.

# COMPLETE DATA MODEL

$$\text{mean of covariates } \boldsymbol{\mu} \sim \text{Normal}\,(0, 100)$$

$$\text{scale of covariates } \boldsymbol{\kappa} \sim \text{Normal}_+\,(0, 2)$$

$$\text{covariates } \mathbf{x}_i \sim \text{Normal}\,\left(\boldsymbol{\mu}, \boldsymbol{\kappa}^2\right)$$

$$\text{regression coefficients } \boldsymbol{\theta} \sim \text{Normal}\,(0, 100)$$

$$\text{observation noise scale } \sigma \sim \text{Normal}_+\,(0, 2)$$

$$\text{outcomes } \mathbf{y} \sim \text{Normal}\,\left(\mathbf{X}\boldsymbol{\theta}, \sigma^2\right)$$

Speaker notes

- We develop a model for every data element, including
  - the outcomes $\mathbf{y}$ in the second block: linear regression conditional on regression coefficients $\boldsymbol{\theta}$, features $\mathbf{X}$, and noise variance $\sigma^2$
  - features $\mathbf{X}$ in the first block: independent normal model conditional on population mean $\boldsymbol{\mu}$, e.g., mean age or blood sugar, and scale $\boldsymbol{\kappa}$, e.g., age or blood sugar variation.
- Beware of change in notation: We previously used $\mathbf{y}$ to denote all data. Here we use $\mathbf{y}$ for outcomes and $\mathbf{X}$ for features.
- We could implement a Gibbs sampler and iterate through all parameters. But we want to apply our Bayesian analysis skills to real-world problem rather than implement Gibbs samplers. Let's use Stan.

```stan
data {
    int n, p;
    vector [n] y;
    matrix [n, p] X_obs;
}

transformed data {
    int n_mis = 0;
    for (i in 1:n) {
        for (j in 1:p) {
            n_mis += is_nan(X_obs[i, j]);
        }
    }
}
```

- As `data`, we declare number of observations `n`, number of features `p`, outcomes `y`, and observed design matrix `X_obs`. Missing elements of `X_obs` are encoded by setting them to `nan`.
- The `transformed data` block is executed once per program to evaluate deterministic transformations of the data. Here, we simply count the number of missing elements `n_mis` in the observed design matrix `X_obs`.
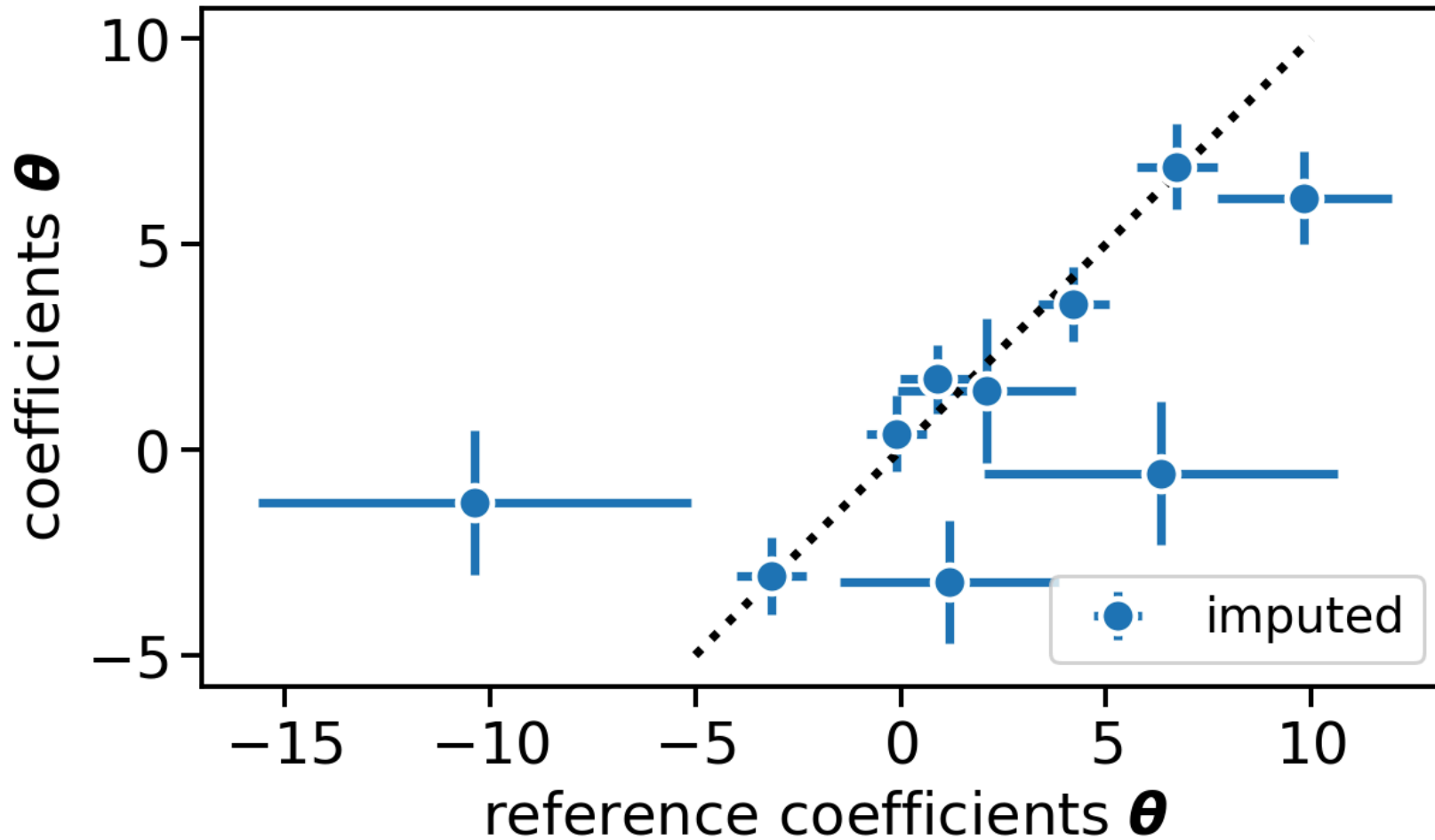
- We declare population mean `mu` and scale `kappa` for each column of the design matrix.
- The regression coefficients `theta` and `sigma` are standard for linear regression models.
- Finally, we have an array of real numbers `X_mis` with size equal to the number of missing elements `n_mis`.

```
parameters {
    vector [p] mu;
    vector<lower=0> [p] kappa;
    vector [p] theta;
    real<lower=0> sigma;
    array [n_mis] real X_mis;
}
```

```stan
transformed parameters {
    matrix [n, p] X;
    {
        int k = 1;
        for (i in 1:n) {
            for (j in 1:p) {
                if (is_nan(X_obs[i, j])) {
                    X[i, j] = X_mis[k];
                    k += 1;
                } else {
                    X[i, j] = X_obs[i, j];
                }
            }
        }
    }
}
```

- In the transformed parameters, we construct the imputed design matrix `X` by copying from `X_obs` if data are available (line #12) or using a latent missing data parameter `X_mis` (#9).
- We use `k` to keep track of the next element we should take from `X_mis`.
- Everything is wrapped in braces `{ ... }` to hide the variable `k` from the `transformed parameters` block. Stan does not allow discrete parameters and complains if `k` is declared outside braces.
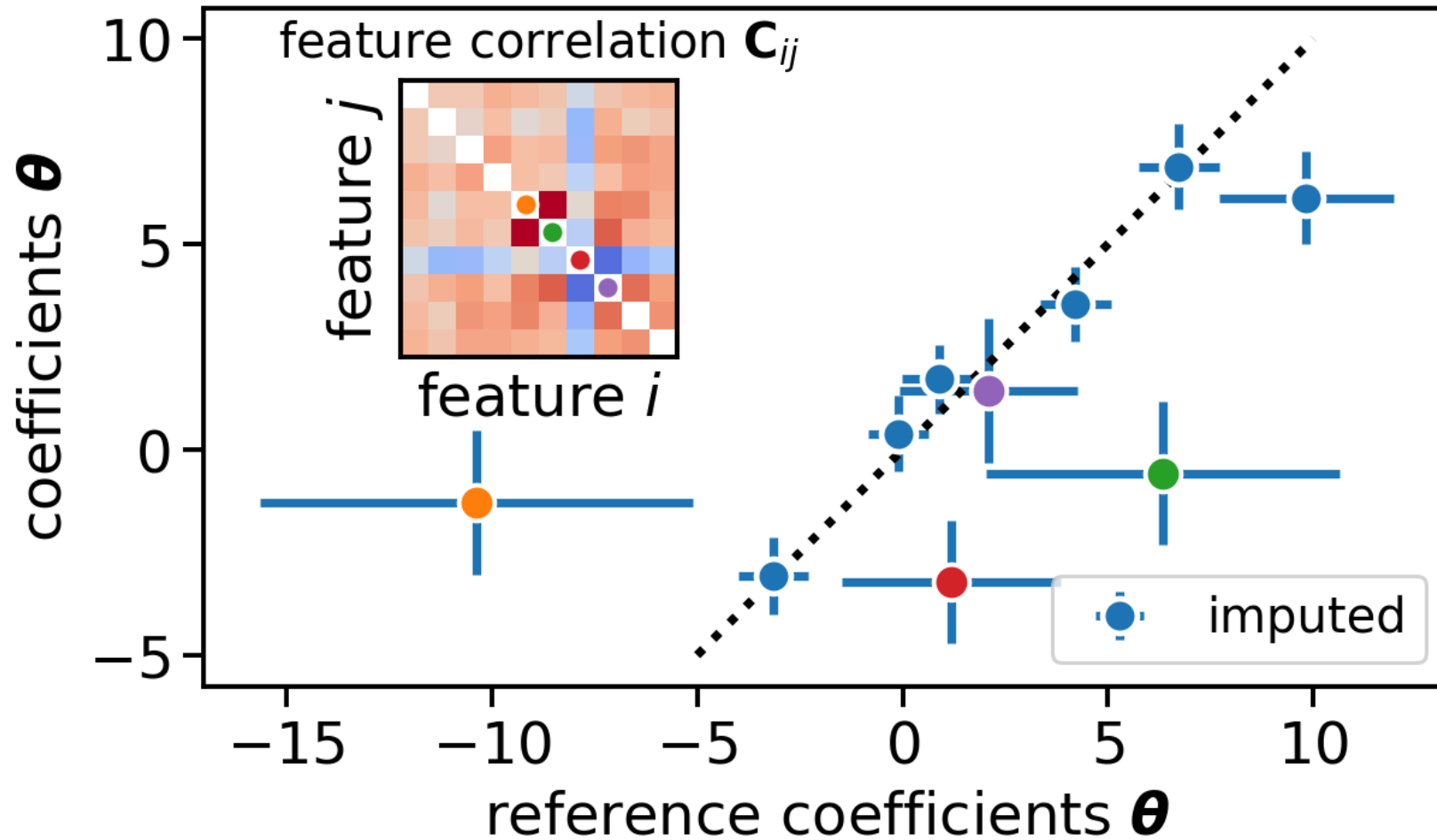
```
1   model {
2       mu ~ normal(0, 100);
3       kappa ~ normal(0, 2);
4       for (i in 1:n) {
5           X[i] ~ normal(mu, kappa);
6       }
7
8       theta ~ normal(0, 100);
9       sigma ~ normal(0, 2);
10      y ~ normal(X * theta, sigma);
11  }
```

- We run the inference twice: (a) on the observed data with elements MCAR, (b) on the complete data to which we have access because we artificially removed elements from the design matrix for this example.
- We plot regression coefficients inferred using observed data against regression coefficients inferred on the complete dataset. Error bars are posterior standard deviations.
- They largely agree, but there are a few disagreements.

- Why is this? Our missing data model is poor: (1) assumed a normal model which may not be appropriate for some features (like sex), (2) features are correlated and an independent normal model for each feature is not appropriate.
- 2nd is problematic: If a person suffers from diabetes, several blood tests are likely abnormal and features are correlated.
- We further investigate by plotting the feature correlation as an inset heatmap. We use colored markers on the diagonal of the correlation matrix to identify features and also color the corresponding regression coefficients.
- We observe that coefficients corresponding to correlated features exhibit disagreements. The issue is likely due to us neglecting feature correlation in the missing data model.

# EXAMPLE: 1988 ELECTION

In 1988, Democrat Dukakis was far ahead of Republican Bush in the polls but decisively lost the election. What happened?

- **Data**: Surveys of public opinion from 51 national polls conducted by nine polling organizations over six months.
- **Goal**: Understand temporal evolution of vote intention for subgroups of the population.
- **Challenge**: Not all questions were asked in all surveys and some questions were not answered.

# BUILDING AN IMPUTATION MODEL

- There are two extremes: complete pooling (impute all surveys together using the same model) and no pooling (impute all surveys independently—not possible for missing questions).
- We use a hierarchical imputation model: imputations are informed by the specific survey when possible and by the population if not.
- This can address both questions not answered (within-survey imputation) and questions not asked (between-survey imputation).

# HIERARCHICAL MODEL

- $y_{siq}$ is response of individual $j$ in survey $s$ to questions $q$.
- $\mu_{sq}$ is the mean response to question $q$ in survey $s$.
- $\mathbf{x}_s$ is a feature vector for survey $s$, e.g., time.
- $\boldsymbol{\beta}$ is a coefficient *matrix*, capturing dependence of responses on features.
- $\boldsymbol{\Psi}$ is a covariance matrix for responses, e.g., opinion of Dukakis may be anticorrelated with opinion of Bush.

$$\mathbf{y}_{si} \sim \text{Normal}\left(\boldsymbol{\mu}_s, \boldsymbol{\Psi}\right)$$
$$\boldsymbol{\mu}_s \sim \text{Normal}\left(\boldsymbol{\beta}\mathbf{x}_s\beta, \sigma^2\mathbf{I}_Q\right)$$

- We have a *matrix* of coefficients because the mean response $\boldsymbol{\mu}_s$ is itself a vector rather than a scalar. The second line is a latent regression model with multivariate outcome.
- We ignore the missing data mechanism because the majority of missingness is due to some surveys not asking certain questions rather than people not answering questions they were asked.

# GIBBS SAMPLER

- Sample missing answers in $\mathbf{y}$.
- Sample answer covariance $\mathbf{\Psi}$.
- Sample latent survey means $\boldsymbol{\mu}$.
- Sample between-survey variance $\sigma^2$.
- Sample regression matrix $\boldsymbol{\beta}$.

# RECAP

- Missing data are just more parameters.
- Distinction between MCAR, MAR, and NMAR.
- Ignorability of the missingness *mechanism* for MAR but *not* the missingness itself.
- Imputation using a Gibbs sampler.
- Stan model for imputed data.