# Methods for Selection Bias in the UK Biobank

Valerie Bradley

University of Oxford, Department of Statistics

15 October 2019

# The UK Biobank

**UK Biobank**

- Largest ever prospective health study
- Includes genetic sequencing, blood tests, physical exams, health history questionnaire
- 500,000 participants aged 40-70 when recruited between 2006 and 2010
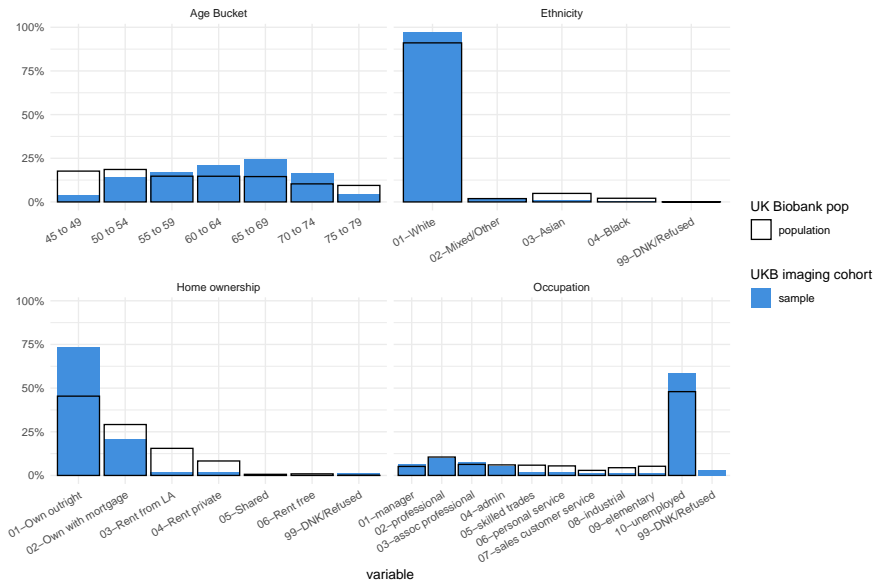
**UK Biobank imaging cohort**

- Subset of UKB participants recruited to undergo additional imaging exams

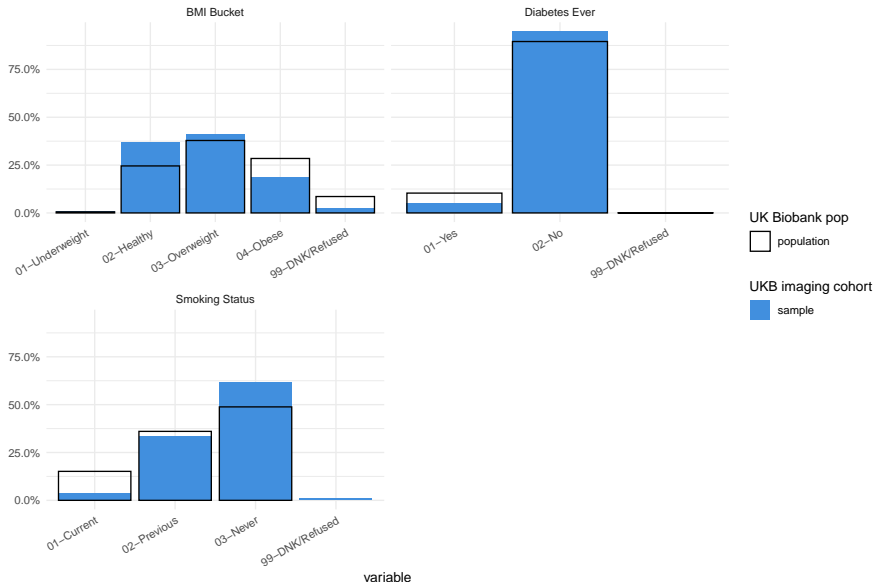**Goal**: Study exposures and outcomes that affect aging populations

# Selection bias in the UK Biobank

## Demographic covariates
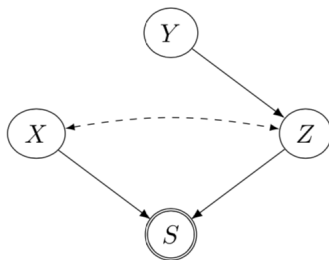
# Selection bias in the UK Biobank

## Health covariates

## Why is this a problem?

For example,

- Y is an individual's hippocampal volume
- X represents a subject's socio-economic status (SES)
- Z represents the subject's level of dementia
- S is selection into the study



By conditioning on selection, if we know that someone is of low SES, they are less likely to show signs of dementia (than if we didn't know their SES)

## Recovering from selection bias

If we can re-write our outcome of interest $P(\mathbf{y}|\mathbf{x})$ as

$$P(\mathbf{y}|\mathbf{x}) = \sum_{\mathbf{z}} P(\mathbf{y}|\mathbf{x}, \mathbf{z}, S=1) P(\mathbf{z} \setminus \mathbf{z^T}|\mathbf{z^T}, S=1) P(\mathbf{z^T})$$

then it's possible to recover from selection bias.

- $\mathbf{z}$ is a set of observed *auxiliary variables*
- $\mathbf{z}^T$ is the subset of $\mathbf{z}$ for which we have external, unbiased population data
- Key is that conditioning on $\mathbf{z}$ makes the outcome of interest, $Y$, *conditionally independent of selection*, $S$

**Problem**: can only condition on a limited number of discrete variables before this breaks down

Methods for selection bias seek to solve this in different ways

# Methods: post-stratification

## Intuition

# Methods: summary

Classic weighting methods:

1. *Post-stratification*: adjust to the joint distribution of **z**
2. *Raking*: iteratively adjust the marginal distributions of elements in **z**
3. *Calibration*: raking, but with continuous variables as well as discrete variables
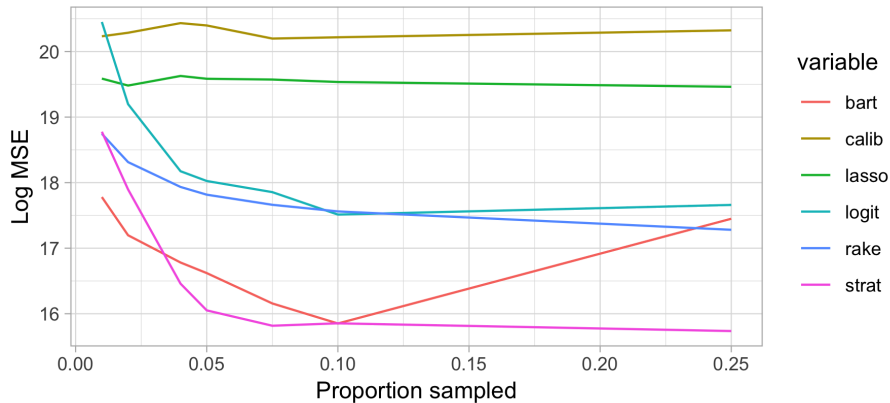
Less-common methods:

4. *LASSO*: use a LASSO to select variables and interactions for raking
5. *Logit*: estimate the probability of selection directly

New method:

6. *BART + raking*: use a BART to estimate the probability of selection, then rake such that key marginal distributions match those of the population
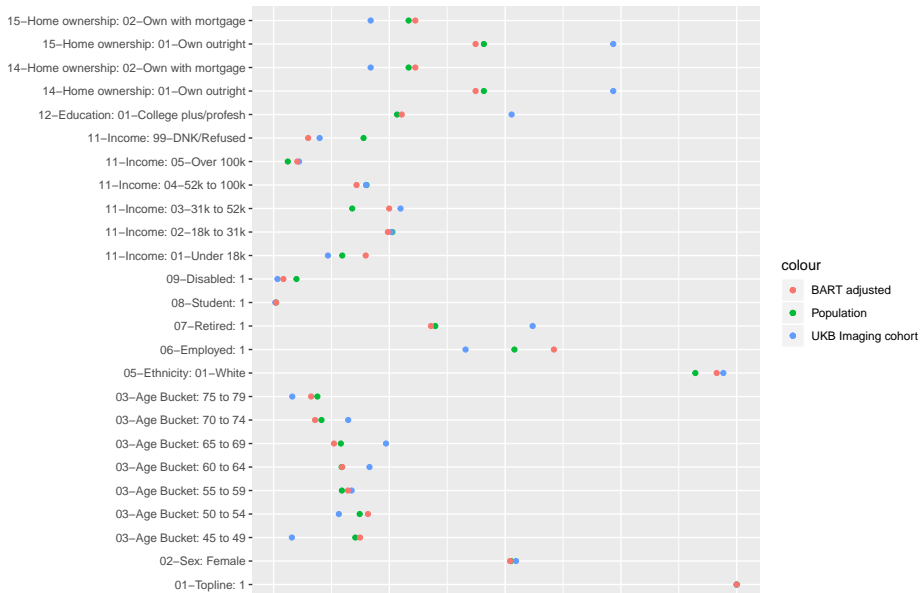
MSE of total brain volume by proportion sampled

UK Population v. UKB Imaging cohort

# Selection bias in the UK Biobank



UK Population v. UKB Imaging cohort