

# **STA 723: Case study 1**

January 17, 2020

*Professor David Dunson*

**Olivier Binette, Joe Mathews and Brian Kunder**

## Exploratory Data Analysis

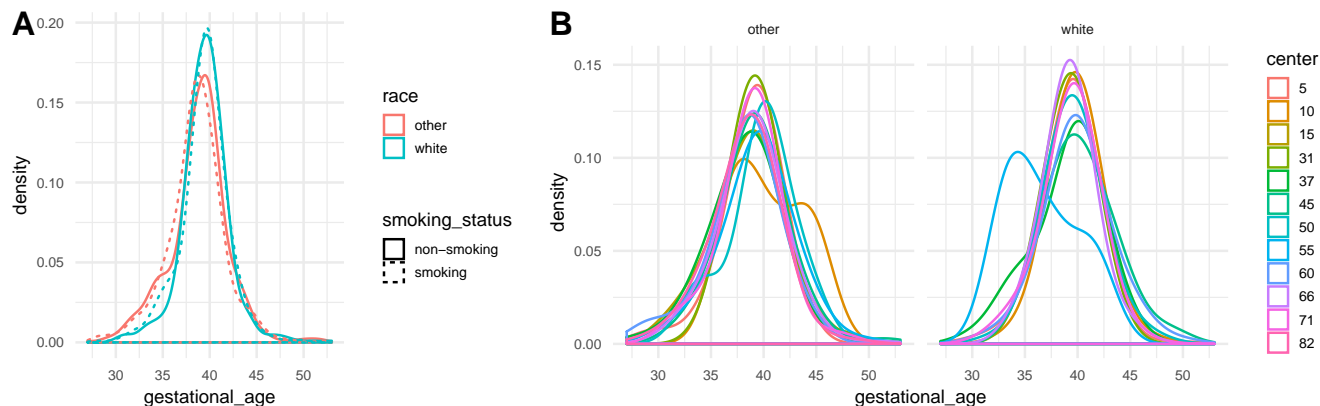
This dataset from the National Collaborative Perinatal Project (CPP) relates gestational age to chemical exposure (DDE and PCBs) and other factors (socio-economic and health-related) in 2380 pregnant women. The goal is to assess how exposure to DDE and PCBs impact the risk of preterm birth, defined as delivery before 37 weeks.

### Data cleaning and manipulations

We removed pregnancies over 55 weeks from the dataset. We also dropped the `albumin` variable, which contains 93% missing values. The variables `score_income`, `score_occupation` and `score_education` contain 21% missing values. Otherwise the data contains only one observation ( $n = 1857$ ) with missing `PCBs` values, which we remove. Given the low sample size for the “other” `race` classification ( $n = 123$ ), we combined “black” and “other” in a single class of size 1336. Finally, to help with visualization and interpretation, we summarized the different PCBs with a positively weighted average. The weights were chosen to minimize the sum of squared orthogonal residuals to normalized PCBs. This ensures that this one-dimensional summary is a relatively good approximation to the PCBs data.

### Gestational age

Gestational age distribution by groups.



Relationship with DDE.