

Residency Exams

Woonggyu Jin

Jack Steel

Sonya Eason

Background & Motivation

Prior to practicing full-time medicine, internal medicine students must complete a three year training known as a medical residency. At the end of their medical residency, they are required to pass an examination to obtain their MD. These residency programs are known for their extensive work weeks and rigorous workload.

In 2003 and 2011, two reform policies were passed to change the structure of internal medicine residency. The 2003 reform, passed on July 1, sought to limit the number of hours worked each week by capping it at 80 hours. The 2011 reform, also passed on July 1, placed stricter limits on which types of activities are permitted during the 80-hour work period.

These reforms were passed with the goal of improving patient care by decreasing the stress placed on internal medicine residents. In other words, with more time away from the hospital to rest and rejuvenate, students will be more successful at their jobs when they are on their shift. However, limiting the number of hours medical residents can work each week brings about several concerns regarding their performance of the examination at the end of their residency. While some believed a work time limit granted them more time to study, others thought that less hands on work would result in decreased pass rates.

These concerns bring about out motivating question: is there empirical evidence of an association between the reforms and the rate at which the medical residents passed the exam?

Data Manipulation and Exploratory Data Analysis

The provided data set contains three columns and twenty rows, containing a column for year, number of exam takers, and pass rate as a decimal. Each row represents a single year in which internal medicine residents completed their residency and took the examination.

We manipulated the data by adding three additional columns for the number of students that passed the examination, the number of students that failed the examination, and the time period based on the year. The first time period is from 1996-2002, when there was no reform policy in place. The second time period is 2003-2010, when the first reform policy was in place. The third time period is 2011-2015, when the second reform policy was in place (*see: Table 1*).

We set time period two as the baseline to make it easier to compare pass rates before and after each policy change. With this setup, the model coefficients directly show whether pass rates were higher in the period before vs. after the first reform, and before vs. after the second reform.

Figure 1 displays the medical resident exam pass rate over the years 1996 to 2015. Because the effects of policy changes might take time to appear, given that internal medicine residency lasts three years, we also considered phase-in periods. To capture this, we divided the plot into five sections. The original three time periods (no reform policy, reform policy one, and reform policy two) were each split further into a phase-in period and a full policy period. The phase-in period represents the three years when residents experienced a mix of the old and new policies during their training.

Students that fall in the section “Phasing in Reform One” (green dots on Figure 1), had experienced both no reform policy and reform policy one. Students that fall in the section “Phasing in Reform Two” (teal dots on Figure 1), had experienced both reform policy one and reform policy two during their residency.

As can be seen from Figure 1, students that took the examination in years where they experienced some amount of reform policy one appeared to perform better than students that took the examination in years where they experience no reform policy or some amount of reform policy two.

Model Implementation Details

Model 1 - Binomial Mixture

We reasonably assumed that exam pass rates contain unobserved year-to-year variation, such as differences in exam difficulty or cohort quality. To capture this extra source of randomness, we fit a binomial mixture model by including a random intercept for each year. This approach allowed us to separate systematic effects of reform policies from idiosyncratic annual noise.

We fit a generalized linear mixed model (GLMM) with a logit link, specifying pass/fail outcomes as the response and reform time period as the fixed effect, with year as a random effect (*see: Equation 1*)

Model 2 - Beta-binomial

The binomial assumption may underestimate variability in exam outcomes because pass probabilities likely vary within each period. To allow for overdispersion, we modeled the yearly pass probabilities as Beta-distributed. This approach treats exam outcomes as binomially distributed, with their underlying pass rate following a Beta distribution, allowing variability beyond what a simple binomial model would capture.

We fit a beta-binomial regression with time period as the predictor using the `glmmTMB` package with a logit link. The model structure is described in equation 2.

Here, dispersion was estimated directly from the data, accounting for extra variation.

Model 3 - Beta-binomial on subset

Because policy reforms were phased in gradually, exam cohorts in transition years likely had mixed exposure. To avoid contamination, we fit the beta-binomial model on a restricted dataset excluding these phase-in years. This design isolates the effect of reforms once they were fully implemented.

We restricted the dataset to years 1996–2002, 2006–2010, and 2014–2015. The same beta-binomial specification was used, with time period as the predictor and a logit link for the mean pass probability.

Model Evaluation

To check goodness of fit, we plotted Pearson residuals over time. We appear to minimize residuals over time in the mixture model, as seen in Figure 2. Putting this together, we concluded that the full beta-binomial mixture model is the best fit (*see: Figure 2*).

Shortcomings

One challenge was finding evaluation metrics that could be compared across models. We first tried using AIC, but this was not reliable since different packages calculate it differently, and AIC depends on the likelihood, which changes when models use different numbers of observations.

For future work, Bayesian models could be helpful since they would let us incorporate additional information into the analysis. We also plan to explore other comparable measures of model fit and parsimony.

Conclusion

We concluded that the binomial mixture model is the best fit. With that model, we found that the first reform increased the odds of passing, while the second reform was associated with a decline in performance. Both of these were confirmed by significant coefficients (*see: Table 2, Figure 3*)

Appendix

Table 1: First three rows of the dataset after adding pass/fail counts and reform time period.

Year	N	Pct	Pass	Fail	timeperiod
1996	6964	0.82	5710	1254	tp1
1997	7173	0.85	6097	1076	tp1
1998	7348	0.84	6172	1176	tp1

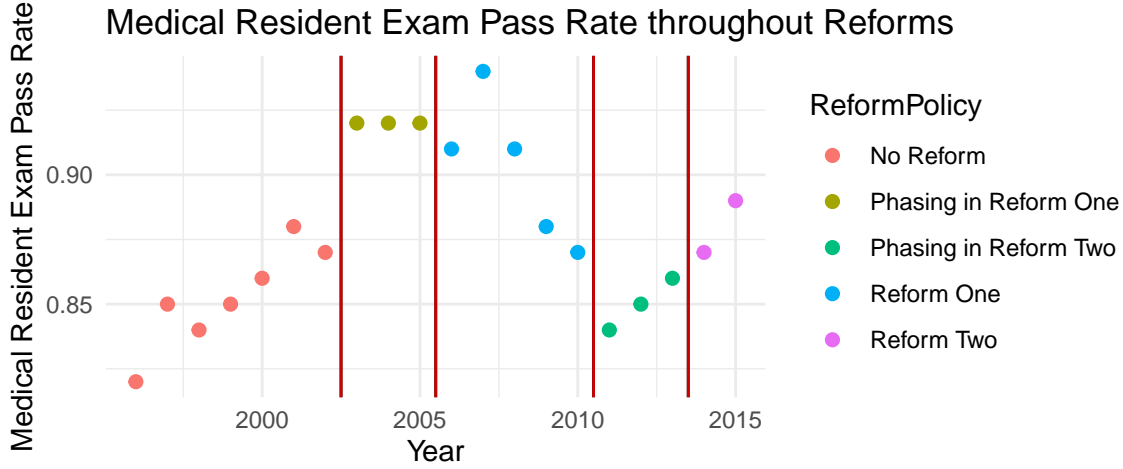
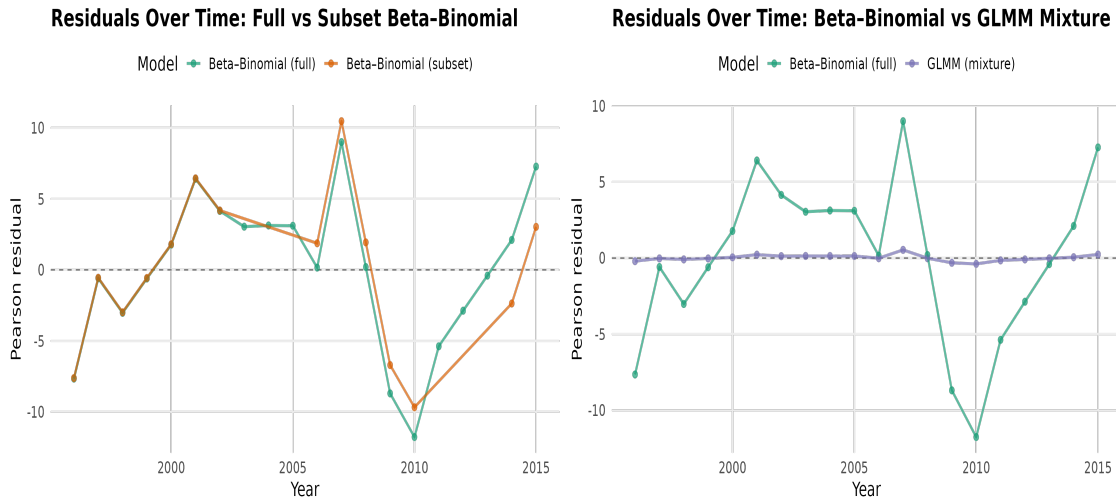


Figure 1: Internal medicine resident exam pass rate over time periods for no reform (1996–2002), phasing in reform one (2003–2005), reform one (2006–2010), phasing in reform two (2011–2013), and reform two (2014–2015).

$$\text{Pass}_{iy} \sim \text{Binomial}(n_{iy}, p_{iy}), \quad \text{logit}(p_{iy}) = \alpha + \beta \cdot \text{timeperiod}_{iy} + u_y, \quad u_y \sim N(0, \sigma^2). \quad (1)$$

$$\pi_y \sim \text{Beta}(\alpha, \beta), \quad \text{Pass}_y \sim \text{Binomial}(n_y, \pi_y). \quad (2)$$



lotted by Year for Full Beta–Binomial Model Compared to Subsetted Model Fits and Mixture Binomial Model

Table 2: GLMM (binomial mixture) fixed-effect coefficients (reference period: tp2).

term	estimate	std.error	statistic	p.value
(Intercept)	2.3234	0.0689	33.7444	0
timeperiodtp1	-0.5591	0.1004	-5.5679	0
timeperiodtp3	-0.4831	0.1103	-4.3786	0

