

# The Determinants of Teachers' Occupational Choice\*

VERY PRELIMINARY

Kevin Lang<sup>†</sup> and Maria D. Palacios<sup>‡</sup>

June 2nd, 2017

## Abstract

Among college graduates, teachers have both low average AFQT and high average risk aversion. Using a dynamic optimization model with unobserved heterogeneity, we find that the low mean AFQT score among teachers does not reflect a low return to cognitive skill within teaching but low return to other skills, correlated with AFQT. The compression of earnings within teaching attracts relatively risk-averse individuals. Were it possible to make teacher compensation mimic the return to skills and riskiness of the non-teaching sector, overall compensation in teaching would significantly increase. Moreover, such a shift would substantially reduce the utility of many current teachers, making the process of reform challenging. Importantly, our conclusions are sensitive to the degree of heterogeneity for which we allow, and even a model with no unobserved heterogeneity appears to fit well within sample. It would be easy to conclude that allowing for two or three types fits the data adequately. Formal methods reject this conclusion. The BIC favors more types through six, the most we have been able to include in the model so far. Ranking models by their out-of-sample accuracy, more types is also better although the improvements from going from four to five and from five to six types are noticeably smaller than those from adding an additional type to a lower base.

---

\*We are grateful to Peter Arcidiacono, Richard Blundell, Ivan Fernandez-Val, Hiro Kaido, Zhongjun Qu and Marc Rysman for helpful discussions and to participants in seminars at Boston University, Carleton University, Queen's University and the University of Western Ontario for their comments and suggestions. Lang acknowledges NSF funding under grant SES-1260917. The usual caveat applies.

<sup>†</sup>Boston University, NBER and IZA. email: lang@bu.edu

<sup>‡</sup>Boston University. email: doloresp@bu.edu.

# 1 Introduction

We make two contributions. On the substantive side, we examine the feasibility of a policy that makes the earnings structure for teachers more closely resemble that outside of teaching. On the methods side, we show that standard tests of in-sample fit can lead us to accept models with inadequate allowance for heterogeneity.

Teacher salaries are typically based on salary scales that depend only on education, experience and seniority and not on measures of quality or achievement. Yet we know that, at least in other settings (e.g. Lazear 2000), tying compensation more closely to performance can both increase the productivity of a fixed set of individuals and attract more productive individuals. Hoxby and Leigh (2004) provide support for this hypothesis in teaching (see also Bacolod (2007)).<sup>1</sup> Consequently, there is considerable interest in performance pay for teachers (see for example, the National Research Council report, Hout and Elliott 2011).

At the same time, teachers unions have typically resisted performance pay. This is not entirely surprising. As we will show, teachers are, on average, more risk averse than the general population of college graduates. Standard theory implies that they will therefore require greater compensation to offset the increased risk associated with performance pay.

In this paper, we do not examine performance pay, *per se*, but we ask how the composition of the teaching profession would change if education and ability were compensated in the same way as in the general labor market for college graduates, presumably making teaching a similarly risky occupation. Like much of this literature, we look at the effect on general ability as measured by test scores and/or potential earnings outside teaching since we do not have a measure of teacher effectiveness. We take it as given that ability within and outside teaching are correlated, albeit imperfectly. This is supported by our estimates of a strong correlation between potential earnings in teaching and non-teaching.

To do this we estimate a dynamic model of occupation choice in which individuals decide each year whether to work as a teacher, in some other occupation, or not to be employed. The decision is dynamic because there are sector-specific returns to experience and because

---

<sup>1</sup>Leigh (2012) finds evidence that higher pay increases the test scores of students in teacher training programs in Australia and some evidence that greater earnings dispersion outside of teaching lowers scores.

there is a cost of moving among the sectors.<sup>2</sup>

Our model is closest in both format and spirit to Stinebrickner (2001b) and particularly Stinebrickner (2001a). There are, however, some important differences in our treatment of uncertainty and of variation in the importance of earnings in the utility function. In addition, unlike him, we do not limit our sample to individuals who obtain their teaching qualification early on but consider all college graduates. The policy changes of interest to us may affect the decision to obtain a teaching qualification. We are also able to follow individuals much later into their careers which allows us to consider exit from and reentry into teaching.

Rewarding this form of ability in teaching similarly to the way it is rewarded outside teaching leads not only to the anticipated shift of teachers to the more skilled types but also raises the AFQT of teachers substantially. Perhaps surprisingly, but consistent with Hanushek, Kain and Rivkin's (2004) finding of only modest effects of earnings on quality, this policy change raises quality by substantially raising average salaries. The effects on the salaries paid to different types is important, but the change actually attracts higher quality workers who are not heavily focused on high earnings.

Even this conclusion ignores the difficulty of effecting a transition. Reform requires transitioning to a compensation system that rewards characteristics differently from the current system and increases risk for a population that is risk averse relative to other college graduates. The reform we study would make a substantial proportion of experienced teachers worse off. Reforms are therefore likely to be very disruptive in the short run, regardless of whether they are beneficial in the long run.<sup>3</sup>

Our conclusions turn out to be very sensitive to the extent to which we allow for unobserved individual heterogeneity. In doing so, we face two risks. If we do not allow for sufficient heterogeneity, the model is misspecified. If we allow for excess heterogeneity, although the estimates remain consistent, our counterfactual estimates may suffer from overfitting of the original model. Strikingly even the model with no unobserved heterogeneity appears to fit the data well. One could easily conclude that allowing for two or three unobserved types

---

<sup>2</sup>There is an enormous literature on the decisions to become a teacher and to leave teaching which we will not attempt to review thoroughly. This literature is reviewed in Dolton (2006).

<sup>3</sup>See Biasi (2016) for a study of teacher mobility and exit following Wisconsin's Act 10, which radically overhauled the compensation system for teachers in that state.

is adequate. But the simulation results are quite different if we allow for four types and dramatically different if we allow for five or six types. To choose the number of types we use both the Bayesian Information Criterion and the accuracy of out-of-sample predictions using different numbers of unobserved types. Both approaches reject using fewer than six types. We have been unable to test whether seven types is preferable to six.

## 2 Data and Some Empirical Regularities

We reverse the usual order of presentation and discuss the data before the model because certain regularities will influence how we develop the model.

We use the National Longitudinal Survey of Youth 1979 (NLSY79). Since the NLSY79 is well-known to labor economists, we skip a general description of the survey. We restrict the sample to college graduates and drop observations for years in which individuals report being self-employed, in the military or working fewer than 35 hours a week. Finally, we drop individuals interviewed fewer than four times. Table 1 shows summary statistics for the 1,085 individuals in our sample, divided into three categories: (1) teachers, (2) non-teachers and (3) not working. Note that an individual might be in all three categories over the course of the panel.

**Compensation over the Lifecycle:** Table 1 shows that teachers, on average, earn less than other college graduates. As shown in figure 1a, teachers' and non-teachers' have similar earnings when 22 years old, but a gap emerges as they age. In this and other figures, we show estimates for individuals age 22-55. However, sample sizes at both extremes are small and should therefore be treated with caution. By the time they are 55 years old, non-teachers earn almost 56 thousand dollars<sup>4</sup> more than teachers do. This pattern does not merely reflect the changing composition of the various groups over the lifecycle. Controlling for individual fixed effects or restricting the estimates to individuals who are only teachers or only non-teachers does not substantially change the pattern.

**Risk and Uncertainty:** The standard deviation of earnings is initially small and similar for teachers and non-teachers, but as age increases the standard deviation remains modest

---

<sup>4</sup>In 2012 dollars.

Table 1: Summary statistics

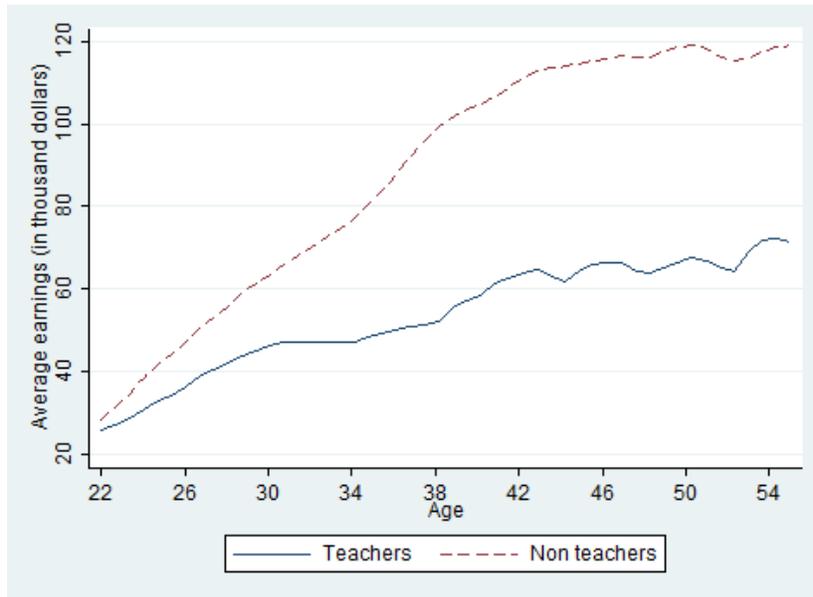
	<b>Teachers</b>	<b>Non-teachers</b>	<b>Not working</b>
	(1)	(2)	(3)
Earnings in \$1,000	51.2387 (0.6978)	76.8924 (0.5519)	- -
Risk aversion standardized	0.2625 (0.0176)	-0.0821 (0.0071)	-0.0197 (0.0350)
Schooling in years	17.4406 (0.0265)	16.8522 (0.0089)	16.7000 (0.0415)
AFQT standardized	-0.3014 (0.0203)	0.0384 (0.0069)	0.0666 (0.0342)
Individuals*	211	1,047	248
Observations	2,619	20,428	810

NOTES: Standard errors in parenthesis.

\*An individual might be in all three categories.

for teachers and grows dramatically for non-teachers (see figure 1b). Since the residual could reflect factors known to the individual but not the econometrician, this does not necessarily imply that teaching is less risky than other occupations, but it is suggestive. And, in fact, if we regress log earnings on schooling, experience, year fixed effects and individual fixed effects, residual earnings variation is higher and grows faster among non-teachers.

(a)



(b)

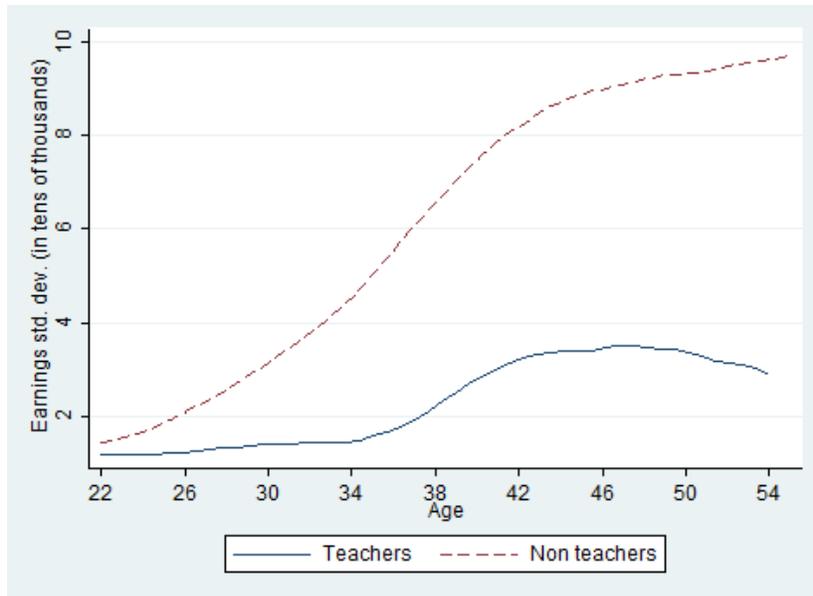


FIGURE 1: Teachers and Non-teachers earnings\* by age

(a) Average earnings

(b) Standard deviation of earnings

*\*Translated into 2012 dollars.*

**Risk aversion:** Therefore, we would expect risk-averse individuals to sort into teaching. We construct a risk-aversion parameter using three questions that were asked in each of four

years:<sup>5</sup> (1) Would you take a job that could double your income or cut it by 1/2 with a 50-50 chance?, (2) Would you take a job that could double your income or cut it by 1/3 with a 50-50 chance?, and (3) Would you take a job that could double your income or cut it by 1/5 with a 50-50 chance? Using the responses to these questions we construct a risk aversion parameter. We assign a “1” to individuals who responded yes to the three questions, then a “2” to individuals who responded no to question one but yes to the other two, then we assign a “3” to individuals who responded no to questions one and two but yes to the last one, and finally we give a “4” to the most risk averse individuals who responded no to all questions. Because the same questions were asked in several years, the risk aversion parameter for a given individual may change. To have a measure of risk aversion for every year we use the most recent risk aversion parameter available for each individual.

As shown in figure 2, teachers are more risk averse than individuals working in other occupations.

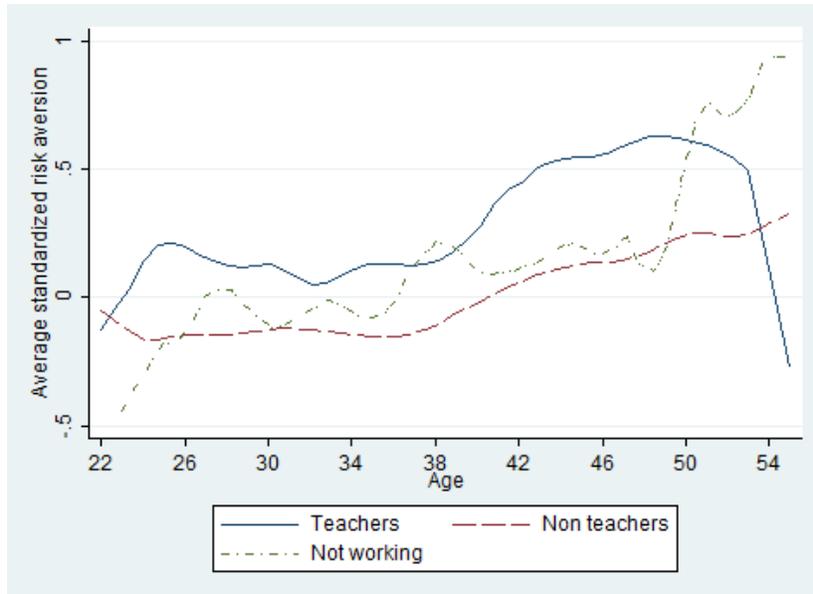


FIGURE 2: Risk aversion by age

**Measured skills:** Since teachers are rewarded for their education, we would not be surprised if they had higher average levels of education than workers in other occupations even among college graduates. This is confirmed in table 1. In fact, in 2000 (when our

<sup>5</sup>1993, 2002, 2004 and 2006

sample is between 35 and 43 years old) 72% of teachers had graduate studies<sup>6</sup> compared with only 43% of non-teachers and 17% of individuals out of the labor force.

However, figure 3 shows that they generally have less skill as measured by the AFQT. The NLSY79 measured the AFQT when the sample was 15-22 years old. The mean AFQT percentile, adjusted for the age at which the individual took that test, for teachers in our sample is 67 while the mean percentile for non-teachers is 75. Figure 3 shows the standardized average AFQT by age. This “observable ability” measure is lower for teachers of all age groups, perhaps because teachers’ earnings are not as responsive to cognitive skills as those for non-teachers. For some age groups this difference is almost 0.5 standard deviations.

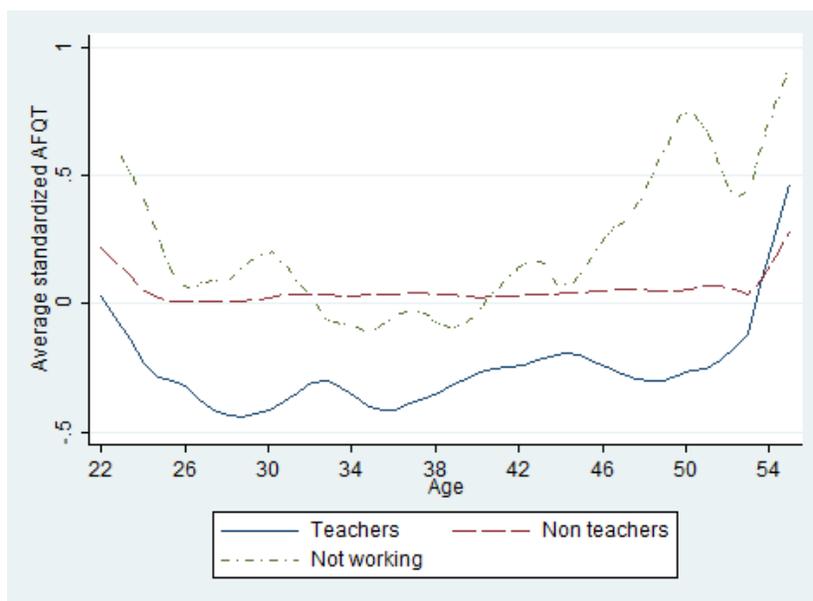


FIGURE 3: AFQT by age

**Occupational mobility:** Table 2 shows that mobility between teaching and non-teaching is modest. Fully 93.4% of individuals who are teachers in one year continue as teachers the following year, while only 5.2% move to a non-teaching occupation even though non-teachers account for 86% of the person-years in our sample. Mobility from non-teaching to teaching is rarer. As can be seen from the second row of the table, only 0.8% of non-teachers transition to teachers in an average year even though teachers account for 11% of the person-years in

<sup>6</sup>Individuals that have at least 17 years of schooling.

Table 2: Transition matrix (percentages)

		$t$		
		Teacher	Non-teacher	Not working
$t - 1$	Teacher	93.4	5.2	1.5
	Non-teacher	0.8	98.3	0.9
	Not working	4.7	23.3	72.0

our sample.

### 3 Model and estimation

Utility is quadratic in earnings:

$$u(w) = aw - bw^2 \quad (1)$$

where  $a > 2bw$ . Letting earnings equal expected earnings plus a mean zero shock with variance  $\sigma_\xi^2$ , expected utility is:

$$E[u] = aE[w] - b(E[w])^2 - b\sigma_\xi^2. \quad (2)$$

Thus, the expected utility of the agent is increasing at a decreasing rate with her expected earnings and decreases with the variance. Although  $b$  does not correspond to a standard coefficient of absolute or relative risk aversion, it captures the worker's degree of risk aversion.

#### 3.1 Occupation choice

We model the occupation choice process as consisting of  $(T - \hat{t})$  periods, where  $T$  is the retirement date, which we will generally impose to be age 60, and  $\hat{t}$  is the first year the individual enters the workforce. Each period, individuals can choose among teaching, other jobs (non-teaching) and non-employment. If they decide to work, their earnings depend on their occupation choice.

Individuals are forward-looking. When choosing an occupation in period  $t$ , they take into account not only the one-period utilities associated with the choices but also their effect on sector-specific experience and therefore on future earnings. In addition, they may face a switching cost.

The one-period utility function depends on the decision made, on expected earnings and on individual characteristics and the switching cost if the individual chooses to switch occupations.

Utility in each state  $d$  depends on the state variables  $s_{it} = \{\bar{Z}_{it}^d, \xi_{it}^d\}$  as given by

$$U_{it}^d(s_{it}) = u(d, \bar{Z}_{it}^d) + \xi_{it}^d \quad (3)$$

$$= \theta_{u1}E[w_{it}^d] + \theta_{u2}(E[w_{it}^d])^2 + \theta_{u3}^d z_{it} + \theta_{u4}(1 - \bar{d}_{i,t-1}) + \theta_{u5}o_{i,t-1} + \xi_{it}^d \quad (4)$$

where  $\bar{Z}_{it}^d := \{E[w_{it}^d], z_{it}, \bar{d}_{i,t-1}, o_{i,t-1}\}$  represents the observed state variables.  $E[w_{it}^d]$  is expected earnings and  $z_{it}$  includes individuals' characteristics such as gender, risk aversion, age, AFQT percentile and schooling.  $\bar{d}_{i,t-1}$  is a dummy variable equal to one if individual  $i$  chose occupation  $d$  or is not employed in period  $t-1$  and zero otherwise. Preliminary investigations found no asymmetry in the cost of switching between teaching and non-teaching. We have therefore imposed that the switching cost  $\theta_4$  is not occupation specific.  $o_{i,t-1}$  is also a dummy variable, equal to one if individual  $i$  chose not to work in period  $t-1$  and zero otherwise. Thus  $\theta_{u5}$  is the cost of transitioning into employment. We normalize the utility associated with non-employment to be 0. We also set the cost of transitioning to non-employment at 0.

Note that individuals' utility for each occupation differs depending on their risk aversion. One should think of the coefficient on risk aversion as corresponding roughly to  $\sigma_\xi^2$  in equation (2).<sup>7</sup> Risk aversion is only measured in four surveys. We hold it fixed between these and assume that individuals know their future risk aversion.

The individual's unobserved preferences for each occupation (taste shocks) are given by the  $\xi_{it}^d$ . We assume the  $\xi_{it}^d$ 's are taken from an extreme value type I distribution. Temporarily assume that the errors are serially uncorrelated. Later, we will address serial correlation by

---

<sup>7</sup>We say roughly because the measure of risk aversion is related to but not identical to the coefficient  $b$  in the expected utility equation.

allowing for multiple types.

From entry into the workforce until retirement ( $T$ ), individuals weigh the consequences of their decisions for future utility. A full solution of the dynamic programming problem consists of finding  $E[\max(V_t^1(s_{it}), V_t^2(s_{it}), V_t^3(s_{it}))]$  at all values of  $z_{it}$ ,  $E(w_{it}^d)$ ,  $\bar{d}_{i,t-1}$  and  $o_{i,t-1}$  for all  $t$ , where the choice-specific value function is:

$$V_t^d(s_{it}) = \begin{cases} U_{it}^d + \delta E[V_{t+1}(s_{i,t+1}) | s_{it}, d_{it} = d] & \text{if } t < T \\ U_{it}^d & \text{if } t = T \end{cases} \quad (5)$$

The choice-specific value function  $V_t^d(s_{it})$  can be decomposed as  $v_{it}^d + \xi_{it}^d$ , where  $v_{it}^d$  is the expected choice-specific value function that has a closed form solution. We set the discount factor  $\delta$  equal to 0.95.<sup>8</sup> Given the extreme value assumption for the distribution of taste shocks, the probability of an individual choosing occupation  $d$  in period  $t$  takes a logit form:

$$P(d_{it} = d | \bar{Z}_{it}, \theta_u) = \frac{\exp(v_{it}^d)}{\sum_d \exp(v_{it}^d)} \quad (6)$$

where the sums are taken over the three possible options available to the individual. Since we have a finite time horizon, and taste shocks are distributed as an extreme value type I, expected value functions have a closed form analytical expression and can be calculated by backward induction. For a more detailed discussion on estimation of discrete choice dynamic programming models see Aguirregabiria and Mira (2010) and Arcidiacono and Ellickson (2011).

### 3.2 Earnings

Earnings depend on the occupation chosen and are a function of: a time trend which also captures linear age effects, individual characteristics, and experience in the teaching and non-teaching sectors. Thus log earnings for a given decision  $d = \{\text{teacher, non-teacher}\}$  in

---

<sup>8</sup>As commented by Aguirregabiria and Mira (2010) the discount factor in most applications is not estimated because it is poorly identified (e.g., see Rust (1987)).

year  $t$  for individual  $i$  are given by:

$$\log(w_{it}^d) = \theta_w^d \bar{X}_{it}^d + \epsilon_{it}^d \quad (7)$$

$$= \theta_{w1}^d f(\text{trend}) + \theta_{w2}^d x_i + \theta_{w3}^d g_{it}^d + \theta_{w4}^d f(\text{exp}T_{it}) + \theta_{w5}^d f(\text{exp}N_{it}) + \epsilon_{it}^d \quad (8)$$

where  $x_i$  includes gender, race and AFQT score.

Schooling,  $g_{it}$ , is occupation specific. For any observed decision, schooling is just the number of years of education that individual  $i$  has, but for the alternative choices we assume individuals would have at least the same schooling as the contemporaneous average individual in that occupation, that is  $g = \max\{\text{actual education}, \text{average education in occupation}\}$ . For instance, if individual  $i$  decides to be a non-teacher in period  $t$  and has sixteen years of schooling then  $g_{it}$  for her non-teaching log earnings equation is sixteen. However, if the average schooling of teachers of her age is higher, say eighteen years, then schooling for individual  $i$  at time  $t$  for her teaching log earnings equation is eighteen. This implies that individual  $i$  would get more schooling if she decided to be a teacher. Since we are not modeling schooling decisions and education is explicitly rewarded in teacher compensation contracts, allowing for education to be higher for off-path decisions is an important feature of a teacher-occupation-choice model.

The occupation-specific experience terms, experience as a teacher ( $\text{exp}T_{it}$ ) and experience as a non-teacher ( $\text{exp}N_{it}$ ), evolve depending on the individual's choices.  $f(\cdot)$  is a quadratic function. Finally, the shocks (the  $\epsilon_{it}^d$ 's) are unknown to the individual at the time of the decision and are assumed to be normally distributed with mean zero and variance  $\sigma_d^2$ . We want to capture the fact that the variance of earnings increases with age (or experience). This is particularly true for the non-teaching occupation. Therefore, we model the  $\sigma_d$ 's as linear functions of age. For identification of the coefficient on earnings in the utility function it is crucial that an exclusion restriction exists; a variable which appears in the earnings equation and only affects utility through earnings. We use the sector-specific experience terms as the exclusion restriction.

### 3.3 Heterogeneity, serial correlation and selection

So far, we have assumed that unobserved preferences and unobserved ability are both uncorrelated over time and uncorrelated with each other. Thus an intense unobservable preference for teaching in period  $t$  would not be related with having an intense unobservable preference for teaching in period  $t + 1$ . Similarly, there is no persistent unobserved ability which is known to the individual but not to the econometrician.

To address these concerns, suppose that there are  $L$  types of people that differ in their preferences for each occupation and in their unobserved abilities.<sup>9</sup> We allow the constant terms of the utility functions and log earnings equations, and the coefficients on the expected earnings terms in the utility functions to vary among types. Thus, the utility and log earnings equations for type  $l$  in occupation  $d$  are:

$$U_{it,l}^d(s_{it}) = \omega_{\mathbf{u},1}^d + \theta_{\mathbf{u}1,1} \mathbf{E}[\mathbf{w}_{it,1}^d] + \theta_{\mathbf{u}2,1} (\mathbf{E}[\mathbf{w}_{it,1}^d])^2 + \theta_{u3}^d z_{it} + \theta_{u4} (1 - \bar{d}_{i,t-1}) \dots \quad (9)$$

$$+ \theta_{u5} o_{i,t-1} + \xi_{it}^d$$

$$\log(w_{it,l}^d) = \omega_{\mathbf{w},1}^d + \theta_w^d \bar{X}_{it}^d + \epsilon_{it}^d. \quad (10)$$

Note that, for a given type, we restrict the coefficients on the expected earnings terms in the utility functions to be the same in teaching and non-teaching. Thus expected earnings in teaching and non-teaching occupations give type  $l$  the same utility.

### 3.4 Estimation

We calculate the parameters using maximum likelihood. Without unobserved heterogeneity, the contribution of individual  $i$  to the likelihood function is the product of the likelihood contribution of occupation decisions  $P(\cdot)$  and the likelihood contribution of earnings  $f_w(\cdot)$ :

$$L_i(\theta) \equiv \prod_{t=1}^{T_i} P(d_{it} | \bar{Z}_{it}, \theta) f_w(\log(w_{it}) | d_{it}, \bar{X}_{it}, \theta). \quad (11)$$

---

<sup>9</sup>See Keane and Wolpin (1997), Eckstein and Wolpin (1999) and Arcidiacono (2004) for other papers that control for unobserved heterogeneity in dynamic discrete choice models. Stinebrickner (2001a) uses this approach in a closely related model of occupational choice by qualified teachers.

To estimate the parameters when we include unobserved heterogeneity we use a mixture distribution, where  $\pi_l$  is the proportion of the  $l$ th type in the population. These proportions and the unobservable preferences and abilities are fixed over time, allowing us to control for serial correlation and selection. With unobserved heterogeneity the contribution of individual  $i$  to the finite mixture of likelihoods is:

$$l_i(\theta, \Omega, \pi) = \log\left(\sum_{l=1}^L L_i(\theta_l, \omega_l) \cdot \pi_l\right). \quad (12)$$

The set of structural parameters to estimate consists of 46 coefficients when there is no unobservable heterogeneity. For each additional type we include there are seven extra parameters. We estimate the model for one, two, three, four, five and six types.

## 4 Results

### 4.1 Choosing the number of types

The most obvious approach to model selection, a likelihood ratio (or similar) test, cannot be used because mixture models violate the requisite regularity conditions because some parameters are not identified under the null. An obvious alternative, Schwarz’s Bayesian Information Criterion (BIC),<sup>10</sup> tends to require a very large number of types, perhaps beyond what is numerically feasible.

In our case, as shown in section 5, visually our model fits the data well regardless of the number of types. We use two formal approaches to choose among the models. First, we calculate the BIC for each specification (see table 3). The BIC continues to improve over the entire range we were able to test.

Second, we use a cross-validation approach. Our approach is based on the following logic. A properly specified maximum likelihood model minimizes out of sample prediction error (Hansen and Dumitrescu, 2016). Therefore, if we believe that one of our models is correctly specified, it should be the one that predicts best out of sample.<sup>11</sup> Thus, we randomly divide

---

<sup>10</sup> $BIC = -2 \cdot \log \text{likelihood} + d \cdot \log(N)$ , where  $N$  is the sample size and  $d$  is the total number of parameters.

<sup>11</sup>This is based on work in progress by Hiro Kaido and Kevin Lang.

Table 3: Bayesian Information Criterion

	BIC
One type	24,912
Two types	19,866
Three types	18,052
Four types	17,426
Five types	17,298
Six types	17,096

our sample into two groups, consisting of 80% and 20% of individuals. We re-estimate the models using the larger sub-sample. Then, we use the new coefficients and the data from the other 20% of individuals and calculate the log-likelihoods for each of the six models. We repeat this exercise twenty times and compare the out-of-sample log-likelihoods.

As with the BIC, this cross-validation approach suggests that we require a large number of unobservable types to address heterogeneity adequately although there is some evidence that we are approaching the requisite number of types. In all twenty replications, four types does better out of sample than three types which in turn fits better than the model with two types which outperforms the model with one type. When we compare six versus five types, the model with more types predicts better out of sample in a clear majority of replications and averaged across the replications, but there are replications that give the opposite result. Interestingly, when comparing five versus four, the average difference in the log-likelihoods of the twenty replications is not statistically different from zero, and only in twelve out of the twenty replications the five type model predicts better out of sample than the four types model. Nevertheless, six types does better than four types in every sub-sample.

Consequently for this early draft, we present the results of the model with six types. For purposes of comparison, we also show the results without unobserved heterogeneity. We continue to explore converging models with a larger number of types.

## 4.2 Estimates of the utility function

The results of the selection equations are given in Table 4.

The first set of rows displays the estimates that are common to all types. Relative to nonemployment, being risk averse increases the utility from teaching and reduces the utility from non-teaching. When we control for unobserved heterogeneity with the six-type model, the point estimates for risk aversion in teachers' and non-teachers' utility functions still suggest that risk averse individuals tend to prefer teaching. In the absence of unobserved heterogeneity, being male appears to increase the preference for teaching over non-teaching. If anything, the opposite holds in the model with six types.

Switching occupations and returning to employment is costly. Controlling for unobserved heterogeneity has very little effect on these coefficients. Using the estimates from the model with no unobserved heterogeneity, for the average individual earning \$74,000, switching occupation is comparable to a decrease of \$10,700 in earnings that year and the cost of returning to employment is equivalent to \$21,600 in that year's earnings.

Age, age squared, AFQT and schooling were also included in the estimation. The coefficients on AFQT are very similar for teachers and non-teachers. Whereas, more educated individuals prefer to work as a non-teacher.

The panels labeled Type x show the relation between types and utility. Each type except the first has an additional utility it receives relative to the first type from each of the occupations. These panels also show the full utility that each type receives from expected earnings (in 10,000's of 2012 dollars) and its square. With only one type, the marginal utility of earnings is increasing up to almost \$200,000, which covers 99% of our observations. When we allow for unobserved heterogeneity, all types except the first and fifth put more weight on earnings at low values than the single type does. The second type only values earnings up to about \$130,000 and the fifth only up to about \$145,000 while the other types have positive marginal utility of earnings up to \$180,000 or higher.

Table 4: Occupation specific utility function parameters

			One type		Six types	
			Coefficient	Stand. Error	Coefficient	Stand. Error
All types	Risk aversion	Teachers	0.0953*	(0.0515)	0.0740	(0.0544)
		Non-teachers	-0.0127	(0.0349)	0.0123	(0.0389)
	Male dummy	Teachers	-0.3207**	(0.1573)	-1.8764***	(0.2252)
		Non-teachers	-0.9451***	(0.1534)	-1.5140***	(0.1960)
	Occupation switching cost		-2.2289***	(0.0613)	-2.1516***	(0.0710)
Cost of returning to employment		-4.4726***	(0.0984)	-4.4291***	(0.1071)	
Type 1	constant	Teachers	27.1235***	(1.9465)	29.8638***	(2.2404)
		Non-teachers	28.8147***	(1.4701)	30.4271***	(1.8702)
	$E[w]$		2.0743***	(0.0949)	1.9810***	(0.1962)
	$E[w]^2$		-0.0559***	(0.0039)	-0.0398***	(0.0083)
Type 2	interaction	Teachers			3.1782**	(1.2576)
		Non-teachers			1.3849	(1.2453)
	$E[w]$			3.1506***	(0.2870)	
	$E[w]^2$			-0.1202***	(0.0225)	
Type 3	interaction	Teachers			7.9904***	(1.9325)
		Non-teachers			5.8235***	(2.0322)
	$E[w]$			1.3189	(0.8241)	
	$E[w]^2$			0.1821*	(0.0953)	
Type 4	interaction	Teachers			0.9941	(1.9208)
		Non-teachers			-8.4416***	(2.5069)
	$E[w]$			2.1159***	(0.2587)	
	$E[w]^2$			-0.0375***	(0.0066)	
Type 5	interaction	Teachers			7.2962***	(1.9164)
		Non-teachers			3.9226	(2.3960)
	$E[w]$			0.8938***	(0.2949)	
	$E[w]^2$			-0.0030	(0.0099)	
Type 6	interaction	Teachers			0.6099	(1.1935)
		Non-teachers			-1.7634	(1.1969)
	$E[w]$			2.8040***	(0.2142)	
	$E[w]^2$			-0.0779***	(0.0110)	

NOTE: Age, age squared, AFQT and schooling were also included.

### 4.3 Estimates of the log earnings equation

Table 5 shows the estimates of the log earnings equations. When we do not consider other forms of heterogeneity, the coefficient on AFQT percentile is lower for teachers. The coefficients are basically identical to each other when we control for unobserved heterogeneity.

Not surprisingly returns to schooling are higher for teachers. As previously discussed, teachers with more years of education receive a higher salary. These coefficients are very similar when we increase the number of types. Also, for both occupation groups males' earnings are higher, but surprisingly the gender earnings gap increases with six types, and is larger for teachers in this specification.

A key to identification of the coefficient on earnings in the utility function is to have a variable which is only in the log earnings regression. We use sector-specific experience as the exclusion restriction. The coefficients on experience basically do not change when including more types. Not surprisingly, teaching experience is particularly relevant for teachers. Teacher earnings increase with teaching experience and continue to do so beyond the range of experience found in our data. Small levels of teaching experience provide little benefit outside of teaching. However, our point estimates suggest that teachers with considerable experience benefit in other jobs.<sup>12</sup> Overall, an extra year of teaching increases yearly earnings around \$1,600 for the average teacher, while it only increases yearly earnings around \$360 for the average non-teacher.

Similarly, experience in other occupations increases earnings for both teachers and non-teachers at a decreasing rate but throughout the relevant experience range non-teaching experience is more valuable outside of teaching.

Finally, the standard deviations are modeled as linear functions of age (i.e.,  $\sigma = \sigma_A + \sigma_B \cdot age$ ). The estimates suggest that the variance of teachers' earnings is lower than the variance of non-teachers' earnings for all age groups.

We comment on earnings differences among the six types in the next subsection.

---

<sup>12</sup>We expect that this reflects very experienced teachers transitioning to other well-paid jobs in education.

Table 5: Log earnings parameters

		One type		Six types	
		Coefficient	Stand. Error	Coefficient	Stand. Error
AFQT percentile	Teachers	-0.0006***	(0.0002)	0.0023***	(0.0003)
	Non-teachers	0.0040***	(0.0001)	0.0022***	(0.0002)
Schooling	Teachers	0.1131***	(0.0038)	0.0994***	(0.0038)
	Non-teachers	0.0731***	(0.0013)	0.0676***	(0.0023)
Male dummy	Teachers	0.1335***	(0.0096)	0.2544***	(0.0117)
	Non-teachers	0.2435***	(0.0033)	0.2222***	(0.0088)
Experience teaching	Teachers	0.0394***	(0.0031)	0.0390***	(0.0025)
	Non-teachers	0.0034**	(0.0015)	0.0039***	(0.0013)
Experience teaching <sup>2</sup>	Teachers	-0.0005***	(0.0001)	-0.0005***	(0.0001)
	Non-teachers	0.0009***	(0.0001)	0.0008***	(0.0001)
Experience non-teaching	Teachers	0.0285***	(0.0016)	0.0312***	(0.0014)
	Non-teachers	0.0873***	(0.0017)	0.0839***	(0.0014)
Experience non-teaching <sup>2</sup>	Teachers	-0.0007***	(0.0001)	-0.0007***	(0.0001)
	Non-teachers	-0.0017***	(0.0001)	-0.0017***	(0.0000)
$\sigma_A$	Teachers	0.2258***	(0.0164)	0.3464***	(0.0128)
	Non-teachers	0.1933***	(0.0086)	0.2692***	(0.0053)
$\sigma_B$	Teachers	0.0031***	(0.0004)	-0.0030***	(0.0003)
	Non-teachers	0.0085***	(0.0003)	0.0017***	(0.0001)
constant	Teachers	8.4372***	(0.0726)	9.0015***	(0.0677)
	Non-teachers	8.6820***	(0.0243)	9.0974***	(0.0389)

Trend, trend squared and race dummies were also included.

Types interactions for the six types model are shown in table 6.

## 4.4 Characteristics of the types

Table 6 displays the average risk aversion and AFQT of the different types, the difference in log earnings equations of each type relative to type 1 in each occupation and the proportions of each type in the teaching and non-teaching groups. The rows are ordered by the size of the type effects in the teachers log earnings equation. Thus type 1 has, *ceteris paribus*, the highest earnings among teachers and type 3 the lowest. The estimates suggest that type 1 has an earnings advantage with respect to other types in the teaching occupation; her teacher earnings are higher than other type's teacher earnings all else equal. Types 1 and 3 are also on average the most risk averse. For the non-teaching occupation, type 4 is the type with the greatest earnings advantage. Moreover, type 4 is the type with the highest average AFQT.

Combining this information, we see that the distribution of types across occupations plays a modest role in the earnings gap between teachers and non-teachers. To see this we can calculate that the average teacher type earns 53 log points less than a type 1 in teaching compared with 18 log points for non-teachers relative to a type 1 in non-teaching. Most of the difference is due to how the types are rewarded. If the distribution of teachers by type were the distribution among non-teachers, their mean earnings would be roughly 6 percent higher. If the type distribution of non-teachers were that of teachers, their mean earnings would be 10 percent lower.

Moreover, there is much more variation in earnings across types outside of teaching than within teaching. The weighted standard deviation of the type effect in teaching is about .28 while it is about .37 outside of teaching. Types that are good at non-teaching also tend to have high earnings in teaching. The correlation is about .74 using the distribution of teachers across types and .93 using the non-teachers' distribution. In short, there is a strong positive correlation between the (to us unmeasured) skills that raise earnings within and outside teaching, but these skills are rewarded more generously outside teaching.

In both teaching and non-teaching the correlation between the type earnings coefficient and the average risk aversion of the type is small although it is somewhat positive for teachers and somewhat negative for non-teachers. More importantly, while we see almost no correlation

Table 6: Average characteristics, interaction coefficients and proportions per type

	Risk aversion	AFQT	Earnings Coefficients		Proportions	
	standardized	standardized	Teachers	Non-teachers	Teachers	Non-teachers
	(1)	(2)	(3)	(4)	(5)	(6)
Type 1	0.0842	-0.0244	-	-	0.1157	0.2221
Type 4	-0.1050	0.3081	-0.2520	0.5164	0.0673	0.0834
Type 5	-0.1105	0.1196	-0.3872	0.2428	0.0733	0.1180
Type 6	-0.0016	-0.0094	-0.4897	-0.2436	0.4055	0.2887
Type 2	-0.0106	-0.1511	-0.7691	-0.5426	0.2816	0.2077
Type 3	0.0876	0.0300	-1.1448	-0.8572	0.0566	0.0801

between a type's average AFQT and its earnings coefficient in teaching, there is a modest positive relation outside teaching. This explains why the differential return to AFQT is noticeably smaller with six-types than in the homogeneous-type model.

## 5 Goodness of fit

It is common in structural papers to examine how well the model matches the data by displaying figures that compare predicted and observed averages. Some papers also calculate in-sample goodness-of-fit statistics. We show that, using these comparisons and statistics, our six models seem to fit most data patterns well: occupation choices, earnings and characteristics of transitioning individuals. Then, we test the accuracy of out-of-sample predictions using the 80-20 division explained in section 4.1. We randomly divide our sample into two groups, consisting of 80% and 20% of individuals. We re-estimate the models using the larger sub-samples. Then, we use the new coefficients and the data from the other 20% of individuals to calculate and test the predictions of our six models. We show that even testing out-of-sample predictions, other than the log-likelihood discussed in section 4.1, no single model clearly outperforms the others.

### 5.1 In-sample fit

We begin by describing the goodness of fit in-sample. Figures 4 and 5 show, for the six-types model, the percentage of teachers and non-teachers by age. Visually the model fits well. These figures are very similar for the six models estimated. Moreover, for every model the  $\chi^2$  goodness-of-fit statistic for each of the choices is below the 5% critical value, also suggesting that the models fit well. There are two important caveats: first it assumes that observations are independent. Yet, given strong persistence in choices, ignoring correlation over time tends to exaggerate the value of the  $\chi^2$  statistic.<sup>13</sup> Second this  $\chi^2$  statistic is not adjusted for the fact that it uses estimated parameters to calculate choice probabilities. Therefore it does not have a chi-squared limiting distribution (Moore, 1977).

---

<sup>13</sup>Consider the extreme case where individuals were always teachers or non-teachers. We would effectively be exaggerating the number of observations by a factor equal to average number of years in the sample.

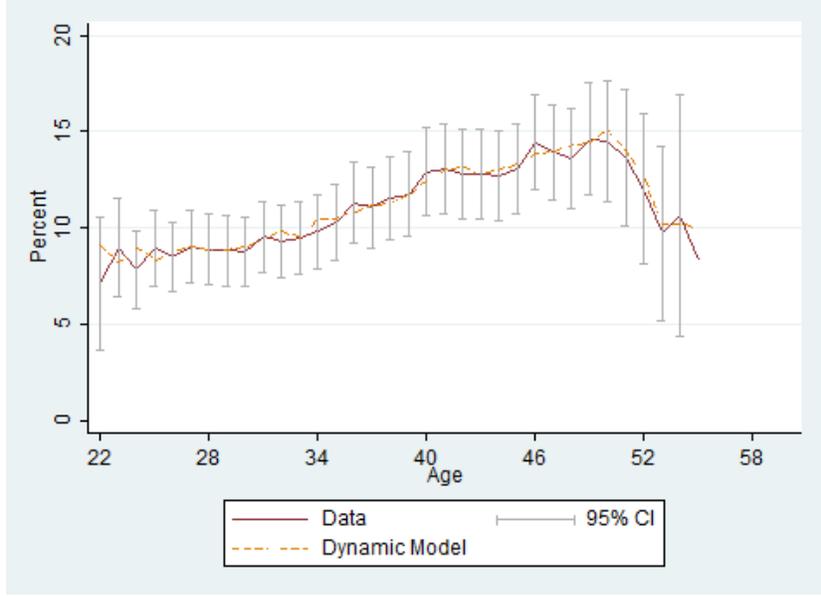


FIGURE 4: Percentage of teachers by age (six-types model)

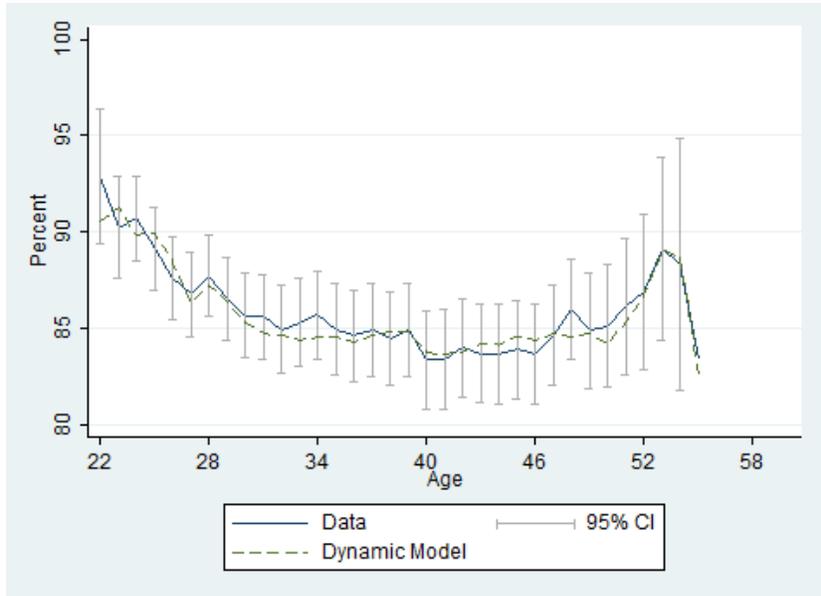


FIGURE 5: Percentage of non-teachers by age (six-types model)

Table 7 shows average earnings for teachers and non-teachers. The model matches the average earnings of both teachers and non-teachers well in the sense that the predicted means lie within the confidence interval of the estimated population means from the data.<sup>14</sup> This is also true for the other five models, the average earnings predicted lie within the confidence interval of the population means. Figure 6 depicts the log earnings by age for both groups,

<sup>14</sup>Standard errors for the model will be calculated for next draft.

Table 7: Average earnings (in \$1,000)

	Teachers	Non-teachers
Data	51.2387 (1.3465)	76.8924 (1.4162)
Model	51.6202	75.5838

NOTE: Standard errors in parenthesis, clustered at the individual level.

for the six-types model. The absolute deviations from the data for different ages go from 0.0195 to 0.3999 for non-teachers and from 0.0416 to 0.4254 for teachers. These numbers are not very different for the other five models.<sup>15</sup>

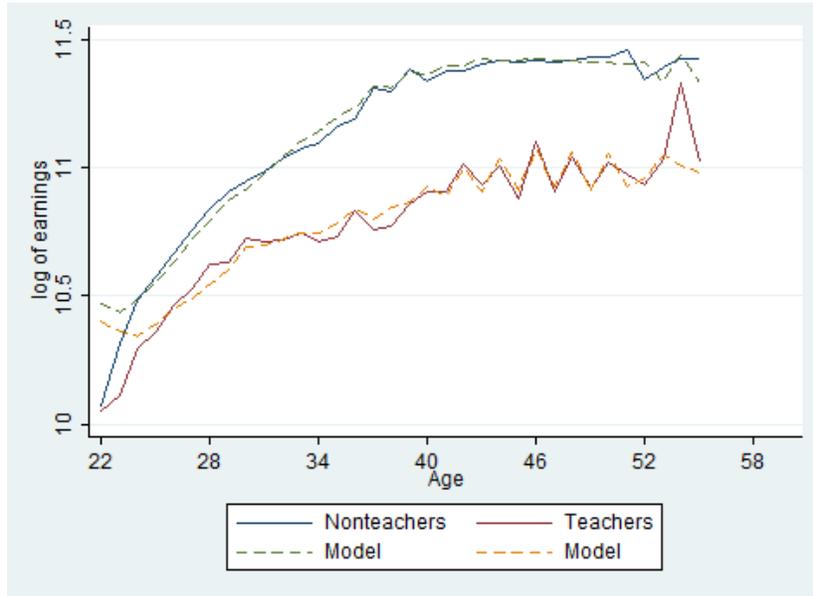


FIGURE 6: Log earnings by age (six-types model)

Comparing tables 2 and 8, we see that the six-types model replicates well transitions for non-teachers and not employed individuals. However, it over-predicts teachers' persistence

<sup>15</sup>The surprising ability to fit the sawtooth pattern at older ages reflects the fact that older workers are only observed in later surveys which were conducted only every other year. Thus the sample for 42 and 44 year olds is roughly constant and disjoint from the one for 43 and 45 year olds. The predictions in each case rely on the observations from in the relevant age group.

Table 8: Transition matrix of predicted choices (percentages)

		predicted choice $t$		
		Teacher	Non-teacher	Not working
choice $t - 1$	Teacher	94.4	3.6	2.0
	Non-teacher	0.9	98.3	0.8
	Not working	5.1	22.7	72.2

and transition into non-employment. Subject to the caveats raised earlier, using a standard goodness of fit test, we reject the model’s fit for transitions for teachers but not for non-teachers and individuals not working. The rejection is largely driven by our under-predicting transitions from teaching to non-teaching employment.<sup>16</sup>

Finally, table 9 shows the average characteristics of transitioning individuals. The model does a good job at matching observed averages in every cell. For every average, the model prediction is within the 95 percent confidence interval for the population average.<sup>17</sup> The same applies for the other five models estimated. An interesting pattern, well replicated by the models, is that individuals who are teachers in one year and continue as teachers the following year are very risk averse and have low cognitive skills. As for non-teachers who do not change occupation, both the observed and predicted averages show that they are not particularly risk averse and have an AFQT score slightly above average.

## 5.2 Out-of-sample fit

While, except for underestimating teachers transitions into non-employment, our models fit well, at least visually, within sample, a fairer test is their ability to predict out-of-sample. We take twenty random samples consisting of eighty percent of the individuals in our sample, reestimate the models and calculate how accurately we predict choices and log earnings for that part of the sample that was not used in the estimation.

When we predict occupation choice, for all five models that allow for any unobservable

<sup>16</sup>Consequently the next draft will drop the assumption that the costs of moving between teaching and non-teaching are independent of the direction of the move.

<sup>17</sup>Using standard errors clustered at the individual level to account for possible correlation of observations.

Table 9: Average standardized risk aversion and AFQT transitions

			<i>t</i>			
			Teacher	Non-teacher	Not working	
<i>t</i> - 1	Teacher	Risk aversion	Data	0.3223	0.1622	0.1415
			Model	0.3350	0.1443	0.1397
		AFQT	Data	-0.3044	-0.4199	-0.3900
			Model	-0.3066	-0.2840	-0.2011
	Non-teacher	Risk aversion	Data	0.1218	-0.0321	-0.1024
			Model	0.1577	-0.0318	-0.0355
		AFQT	Data	-0.3847	0.0408	0.0569
			Model	-0.2056	0.0380	0.0903
	Not working	Risk aversion	Data	0.3072	-0.2140	0.0711
			Model	0.1935	-0.0691	0.0568
		AFQT	Data	-0.2896	0.0720	0.0745
			Model	-0.2954	0.1329	0.0674

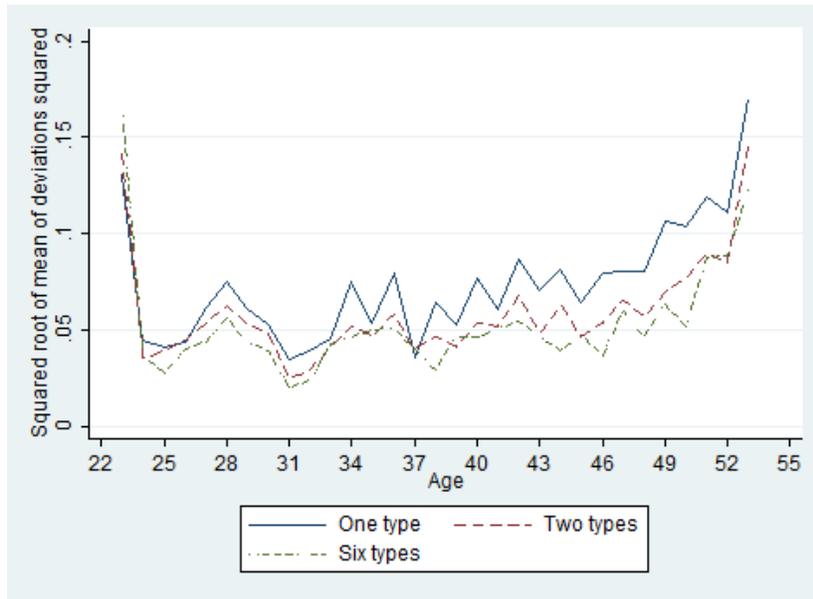
NOTE: \* means the population average is statistically different from the point estimate using SE clustered at the individual level.

heterogeneity, we find that in twenty out of the twenty sub-samples the  $\chi^2$  goodness-of-fit statistic is less than the critical value at a 5% confidence level. Of course, the two caveats raised earlier apply; the  $\chi^2$  statistic does not consider the possible correlation of observations or the fact that we are using estimated parameters to calculate choice probabilities. Nevertheless, this approach is not helpful in distinguishing among the six models.

To assess how well our models predict earnings we follow the spirit of table 7; we calculate average earnings for both groups, teachers and non-teachers, in every sub-sample and compare these averages with the out-of-sample predictions of our models. In the model with no unobservable heterogeneity, in four of the twenty sub-samples for at least one of the two occupations, we reject the hypothesis that the predicted average is within the confidence interval of the observed population average of that group. For the two and three types models, for only one of the twenty sub-samples do we reject the hypothesis that predicted teachers' average earnings is within the confidence interval of the observed population average. For every other model, with four, five or six types, we cannot reject this hypothesis in any of the twenty sub-samples.

Additionally, we analyze the deviations of predicted from observed earnings by age for the twenty sub-samples. We summarize this information in figure 7. These are unconventional graphs; they show the root mean squared error (square root of the mean of deviations squared for the twenty sub-samples) from the one, two and six-types models. There is little difference among the three, four, five and six-types models so we present only the last of these. We also omit ages below 23 and above 53 years where deviations are larger due to a small number of observations. The “average deviations” depicted by age show a moderate decrease when comparing the one type model and either of the models with unobservable heterogeneity.

(a)



(b)

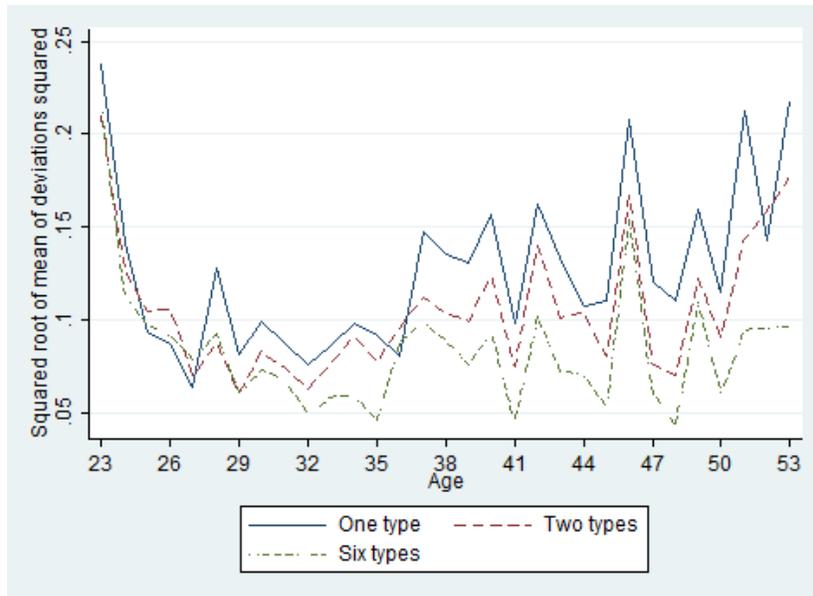


FIGURE 7: Root mean squared error

(a) Non-teachers

(b) Teachers

Finally, we apply our approach to predicting the characteristics of transitioning individuals. In each case, we ask whether the prediction from the model matches the population mean allowing for the confidence interval around the estimate of the population mean but not (in the current draft) the imprecision of the model estimate. We calculate the t-statistic for the

hypothesis that the population average equals the model point estimate.<sup>18</sup> Table 10 shows the maximum t-statistic of the twenty draws for every cell for the six-types model. Since we are performing multiple tests, we use Bonferroni's correction. If there is no correlation among the tests, the adjusted critical value at the 5% confidence level is 3.0233. Since the tests are calculated with sub-samples drawn from the same master sample there should be some correlation among tests, so we can use 3.0233 as an upper bound and 1.96 as a lower bound. That is, we definitely reject the hypothesis that model and data averages are the same for cells with a t-statistic higher than 3.0233, and we definitely cannot reject (at the .05 level) the hypothesis that model and data averages are statistically equal for cells with a t-statistic lower than 1.96. In several cells we definitely do not reject that predicted and observed averages are equivalent; these cells are on the diagonal and for individuals transitioning from teacher to non-employment and from non-teacher to teacher. On the other hand we can only clearly reject the null for average AFQT for the small group transitioning between non-employment and teaching.

Interestingly, the models one, two or three types models predict the characteristics of transitioning individuals slightly better than those with four, five or sixtypes in the sense that there are more cells where averages are definitely not statistically different while in all cases there is only one average for which we can clearly reject the null.

---

<sup>18</sup>Using standard errors clustered at the individual level.

Table 10: Maximum  $t$ -statistic of the 20 draws

			$t$		
			Teacher	Non-teacher	Not working
$t - 1$	Teacher	Risk aversion	0.2592	1.7147	1.7730
		AFQT	0.2502	2.0482*	1.8135
	Non-teacher	Risk aversion	1.4679	0.1230	2.1634*
		AFQT	1.7605	0.1467	1.5685
	Not working	Risk aversion	1.9782*	2.0545*	1.0840
		AFQT	10.8187**	2.1658*	1.0690

NOTE: \*\* denotes cells where averages are statistically different, \* denotes cells where averages could be statistically different.

## 6 Simulation

An advantage of structural modeling over reduced form estimation is the ability to perform counterfactual experiments. However, the credibility of such simulations relies on the model being correctly specified. When allowing for unobserved individual heterogeneity there are two risks. If we do not allow for sufficient heterogeneity, the model is misspecified. If we allow for excess heterogeneity, although the estimates remain consistent, our counterfactual estimates may suffer from overfitting the original model. Therefore, choosing the correct number of types is crucial, particularly if the results from simulations differ depending on the number of types.

In this section we show that the conclusions from our counterfactual experiment are highly sensitive to the number of types. Using our six specifications we simulate how teachers' characteristics, in particular AFQT and risk aversion, would vary if a different contract were offered. To assess the experiment's 'cost' we also calculate average earnings given the new environment. We are interested in analyzing the change in teacher composition keeping quantity constant, that is we want to keep average probabilities of being a teacher and a non-teacher unchanged. To achieve this we adjust the constant terms in the log earning equations for each simulation.

To be clear, we are not simulating an ideal policy. Such a policy would require objective

and/or subjective measures of teaching performance that are not available in our data for teachers, let alone non-teachers. Instead we rely on the positive relation between teaching and other skills to justify this experiment.

Our simulation considers a world in which salaries for teachers reward observable and unobservable abilities in the same way as those of non-teachers. These changes imply an increase in the riskiness of the teaching occupation; there is much more variation in earnings across types outside of teaching than within teaching. To replicate this scenario we adjust several coefficients. First, we adjust teacher salaries so that they are as responsive to AFQT (our measure of “observable ability”) and to schooling as those for non-teachers. Then, we replace the coefficients on the unobservable types in the log earnings equation for teachers with those from the non-teaching log earnings equation. Finally, we set the coefficient on risk aversion in the teacher utility equation equal to its value in the utility equation for non-teaching. We also adjust the constant term so that the utility from teaching for individuals with the lowest risk aversion measure is unaffected by the increased riskiness.<sup>19</sup> We also raise teachers’ earnings variance by setting the  $\sigma_A$  and  $\sigma_B$  parameters equal to the parameters estimated in the non-teaching earnings equation. Due to the conversion from logs to levels, this increases expected earnings in the selection equations.

In summary, we are changing: the risk aversion parameter in the teacher utility equation and the AFQT, schooling, variance coefficients and types coefficients in the teacher log earnings equation. Finally, we adjust the constant terms from log earnings equations so that there is only a change in the composition of teachers and not in the quantity (or average probability).

Table 11 shows the average effect of the simulation using the six models. The results differ dramatically depending on the number of types used, and the effects are not necessarily monotonic in the number of types.

Average AFQT increases only 0.03 standard deviations using the three types model and increases 0.40 standard deviations using the five types model. Figure 8 shows the AFQT change by age under the simulated scenarios. The differences across simulations are partic-

---

<sup>19</sup>We make this adjustment because we would not anticipate that the expected utility of risk-neutral individuals to be affected by the increased riskiness. Given the phrasing of the question, we know only that this group has a low degree of risk aversion and not that its members are risk neutral.

ularly large between ages 30 and 45.

Risk aversion, in turn, decreases by around 0.02 standard deviations with the one, two, and three types models but decreases by more than 0.2 standard deviations with the five and six types models. Figure 9 shows differences in risk aversion by age under the different simulations. Changes are particularly large among teachers in their thirties.

The changes in age, gender and race are negligible for the one through four types models. With five types, the change in average age is nontrivial. It is particularly striking that the share of males increases dramatically in the simulation using the five and six types models. And we see a very notable decline in minority representation with five types and, to a lesser degree, with six types.

How much would such a change in policy cost? Here, again, the answer depends in important ways on the amount of heterogeneity we allow. The smallest ‘cost’ occurs when there are three types, in which case average teaching salaries increase by \$1,500 per year or 2.9%. The estimated costs are also modest with one type (4.0%) and grow somewhat large with two types (7.3%) or four types (9.4%). However, our conclusions are strikingly different with five (27.2%) or six types (24.5%).

What accounts for these differences? Once we have more than one type, a large part of the change in the structure of earnings comes from the way that unobservable types are rewarded. As we include more types, the magnitude of the largest earnings gap for the type with the largest gap will tend to increase. This increases the value of switching from nonteaching to teaching for the group that benefits most from the change. Of course, this is not always the case, and the effect is partially offset by the tendency for the proportion of individuals in each type to fall. Thus the effect need not be monotonic. In our estimates and simulations with five and six types, we see a large shift of the most highly paid group into teaching.

Consistent with this explanation, the disruption required to effect such a policy is much greater in the simulations with five or six types. With one to three types we estimate that 11% - 12% of teachers (measured by teaching years) would leave teaching. When juxtaposed with an annual turnover rate of about 7%, this strikes us as manageable if the policy were phased in over an extended period. The level of turnover becomes somewhat more problematic

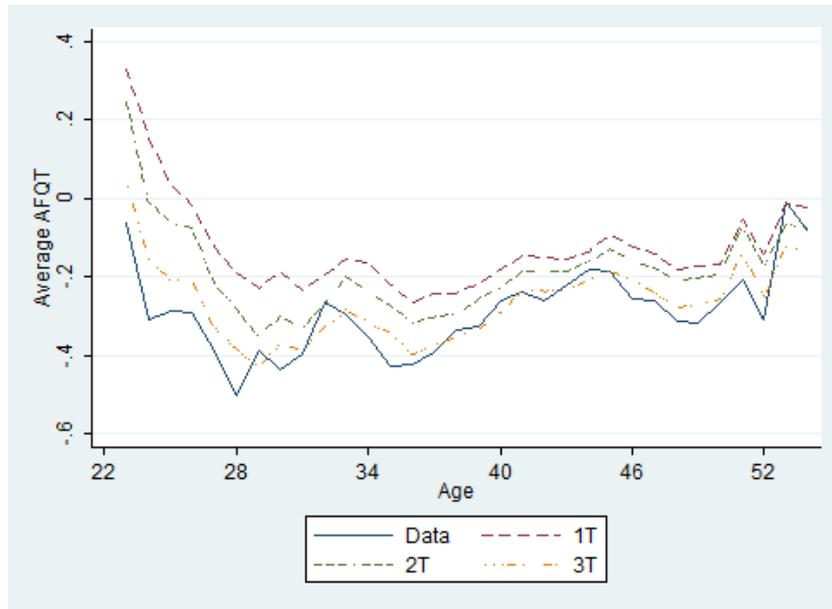
(17%) in the simulation with four types.

While the degree of disruption varies dramatically among the simulations, all suggest to varying degrees significant teacher resistance. In the homogeneous model, 37% of teachers (again weighted by years teaching) would be made worse off. With two types this grows to a majority, and with six types reaches 80%. Even in settings where teachers are not unionized, this would make a transition to this policy difficult.

Table 11: Average effect of simulations

		AFQT	Risk av.	Age	Male	Black	Hispanic	Earnings
		(standardized)						
		(in \$1,000)						
	DATA	-0.3014	0.2625	37.8870	0.2551	0.1672	0.1367	51.2384
One type	Model	-0.2985	0.2610	37.8529	0.2579	0.1651	0.1390	51.5299
	Simulation	-0.1366	0.2376	37.8212	0.2606	0.1335	0.1333	53.5724
Two types	Model	-0.2985	0.2609	37.8490	0.2576	0.1642	0.1397	51.4246
	Simulation	-0.1960	0.2403	37.6179	0.2523	0.1470	0.1356	55.1862
Three types	Model	-0.2987	0.2599	37.8544	0.2578	0.1658	0.1382	51.3414
	Simulation	-0.2727	0.2422	37.5137	0.2562	0.1571	0.1382	52.8461
Four types	Model	-0.3015	0.2605	37.8593	0.2564	0.1676	0.1379	51.3901
	Simulation	-0.2283	0.1910	36.9889	0.2668	0.1519	0.1273	56.2044
Five types	Model	-0.3030	0.2608	37.8768	0.2579	0.1657	0.1361	51.6839
	Simulation	0.0988	-0.0060	35.8158	0.5040	0.1066	0.0866	65.7617
Six types	Model	-0.2988	0.2603	37.8731	0.2589	0.1675	0.1349	51.6202
	Simulation	-0.0428	0.0622	36.4014	0.4334	0.1472	0.0985	64.2852

(a)



(b)

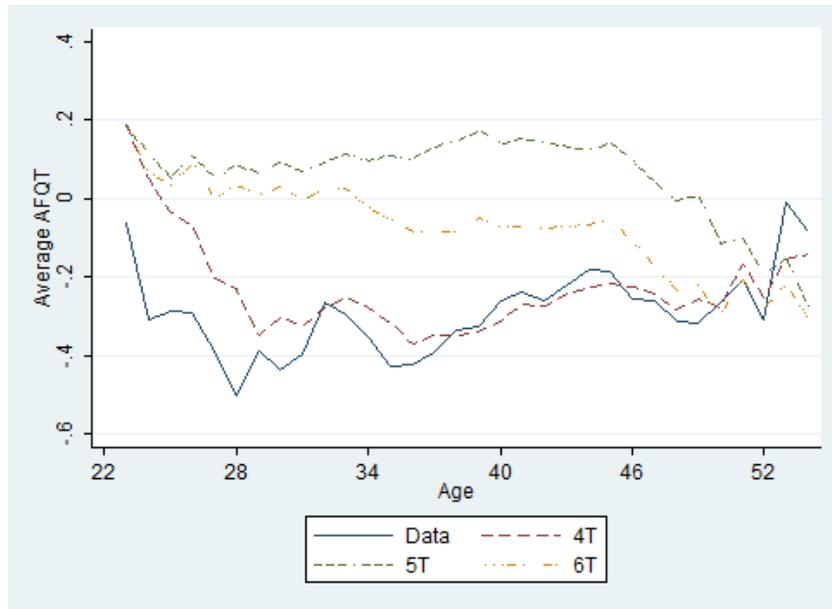
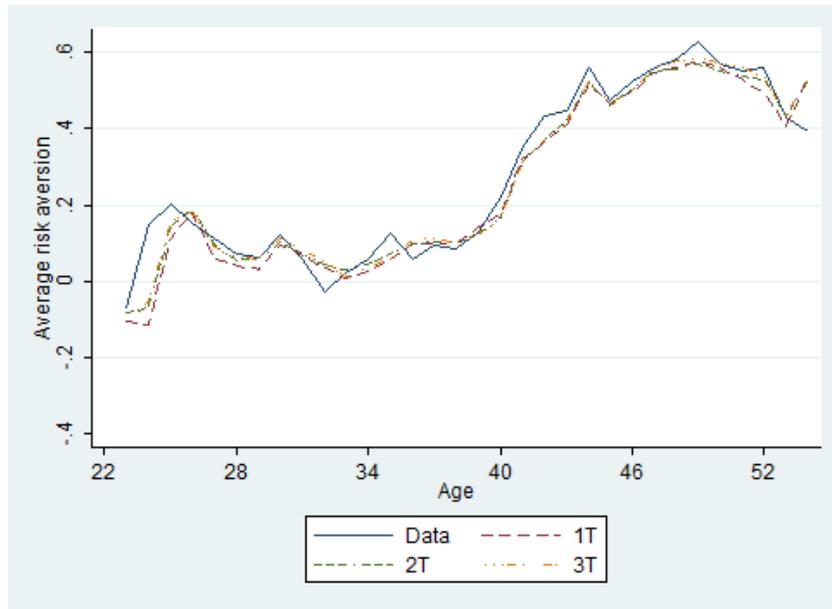


FIGURE 8: Results from simulation on standardized AFQT by age

(a) Using one, two and three types models

(b) Using four, five and six types models

(a)



(b)

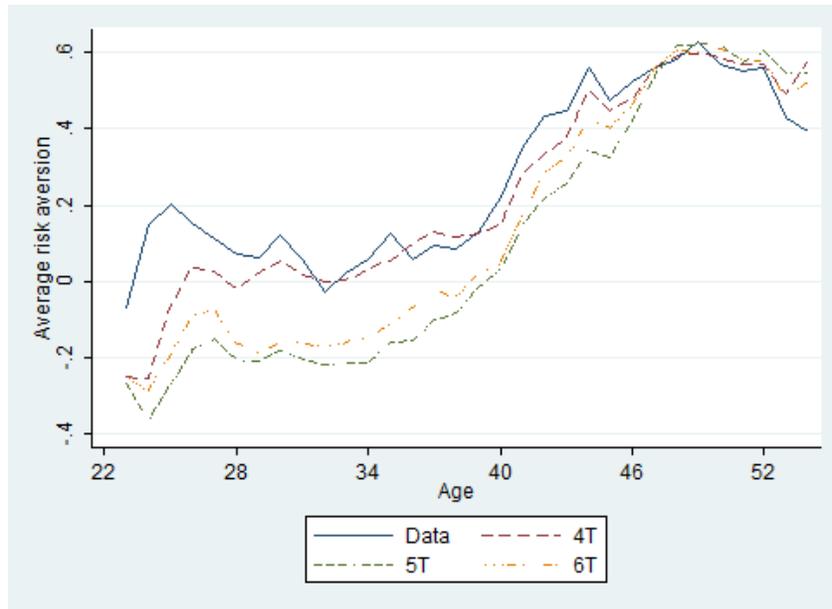


FIGURE 9: Results from simulation on standardized risk aversion by age

(a) Using one, two and three types models

(b) Using four, five and six types models

## 7 Conclusion

This paper contributes both to our substantive understanding of reforming teacher compensation and to the practice of structural modeling.

With respect to the latter, it is widely recognized that the strength of structural modeling is that it allows us to consider experiments that lie outside the data. It is equally widely recognized that the validity of the experiment relies on the model being (at least approximately) correct. We show that establishing that a model fits well within sample is, at best, weak evidence that it is approximately correct. Even showing that it fits well out-of-sample should not be convincing. In our case, even a model with no unobserved heterogeneity appears to fit well within sample and fits reasonably well out of sample. It would be easy to conclude that allowing for two or three types is adequate to fit the data. But more formal methods do not support this conclusion. The Bayesian Information Criterion favors more types through six, which is the most that we have been able to include in the model so far. Perhaps more importantly, when we rank models in terms of their out-of-sample predictions, we also conclude that more types is better. Thus, at least in settings such as ours, it is possible to measure the relative out-of-sample performance of multiple models and choose the one that does best. As we show, the conclusions we draw from our experiment can depend crucially on choosing among models, all of which appear to fit well both within and out-of-sample.

With respect to reforming teacher compensation, we establish that, among college graduates, teachers are not only drawn disproportionately from the lower part of the AFQT distribution, but they are also more risk-averse than their counterparts outside teaching. When we allow for unobserved heterogeneity, the low mean AFQT score among teachers reflects not a low return to cognitive skill within teaching but low returns to other skills, correlated with AFQT. The compression of earnings within teaching attracts relatively risk-averse individuals.

Subject to the very important caveat that we have been unable to go beyond six unobservable types, we show that if it were possible to revise compensation in teaching to mimic the return to skills and riskiness of the non-teaching sector, the effect on overall compensation

in teaching would be considerable. Another difficulty is that such a shift would adversely affect many of those who are currently in teaching and who would suffer large utility losses if they shifted out of teaching. This makes the process of reform challenging.

## References

- AGUIRREGABIRIA, V., AND P. MIRA (2010): “Dynamic discrete choice structural models: A survey,” *Journal of Econometrics*, 156(1), 38–67.
- ARCIDIACONO, P. (2004): “Ability sorting and the returns to college major,” *Journal of Econometrics*, 121(1), 343–375.
- ARCIDIACONO, P., AND P. B. ELLICKSON (2011): “Practical methods for estimation of dynamic discrete choice models,” *Annu. Rev. Econ.*, 3(1), 363–394.
- BACOLOD, M. P. (2007): “Do alternative opportunities matter? The role of female labor markets in the decline of teacher quality,” *The Review of Economics and Statistics*, 89(4), 737–751.
- BIASI, B. (2016): “Unions, Salaries, and the Market for Teachers: Evidence from Wisconsin,” .
- DOLTON, P. J. (2006): “Teacher supply,” *Handbook of the Economics of Education*, 2, 1079–1161.
- ECKSTEIN, Z., AND K. I. WOLPIN (1999): “Why youths drop out of high school: The impact of preferences, opportunities, and abilities,” *Econometrica*, 67(6), 1295–1339.
- HANSEN, P. R., AND E.-I. DUMITRESCU (2016): “Parameter Estimation with Out-of-Sample Objective, unpublished,” .
- HOUT, M., AND E. ELLIOTT, STUART (2011): *Incentives and test-based accountability in education*. National Academies Press.
- HOXBY, C. M., AND A. LEIGH (2004): “Pulled away or pushed out? Explaining the decline of teacher aptitude in the United States,” *The American Economic Review*, 94(2), 236–240.
- KEANE, M. P., AND K. I. WOLPIN (1997): “The career decisions of young men,” *Journal of political Economy*, 105(3), 473–522.

- LAZEAR, E. P. (2000): "Performance pay and productivity," *American Economic Review*, 90, 1346–1361.
- LEIGH, A. (2012): "Teacher pay and teacher aptitude," *Economics of Education Review*, 31(3), 41–53.
- MOORE, D. S. (1977): "Generalized inverses, Wald's method, and the construction of chi-squared tests of fit," *Journal of the American Statistical Association*, 72(357), 131–137.
- RUST, J. (1987): "Optimal replacement of GMC bus engines: An empirical model of Harold Zurcher," *Econometrica: Journal of the Econometric Society*, pp. 999–1033.
- STINEBRICKNER, T. R. (2001a): "Compensation policies and teacher decisions," *International Economic Review*, 42(3), 751–780.
- (2001b): "A dynamic model of teacher labor supply," *Journal of Labor Economics*, 19(1), 196–230.