

Income and Education of the States  
of the United States: 1840–2000

Scott Baier, Sean Mulholland,  
Chad Turner, and Robert Tamura

Working Paper 2004-31  
November 2004

## Income and Education of the States of the United States: 1840–2000

Scott Baier, Sean Mulholland, Chad Turner, and Robert Tamura

Working Paper 2004-31  
November 2004

**Abstract:** This article introduces original annual average years of schooling measures for each state from 1840 to 2000. The paper also combines original data on real state per-worker output with existing data to provide a more comprehensive series of real state output per worker from 1840 to 2000. These data show that the New England, Middle Atlantic, Pacific, East North Central, and West North Central regions have been educational leaders during the entire time period. In contrast, the South Atlantic, East South Central, and West South Central regions have been educational laggards. The Mountain region behaves differently than either of the aforementioned groups. Using their estimates of average years of schooling and average years of experience in the labor force, the authors estimate aggregate Mincerian earnings regressions. Their estimates indicate that a year of schooling increased output by between 8 percent and 12 percent, with a point estimate close to 10 percent. These estimates are in line with the body of evidence from the labor literature.

JEL classification: O40, J24, E01, N00

Key words: state human capital, state output per worker, returns to schooling

---

The authors thank the workshop participants at Clemson University, Texas A & M, the University of Kentucky, SUNY Buffalo, the University of Virginia, the University of South Carolina, the 2004 Midwest Macroeconomic Meetings, Iowa State University, and joint seminars at UNC-Chapel Hill and Duke University for helpful comments. They benefited from suggestions by Kevin Murphy and Casey Mulligan. The views expressed here are the authors' and not necessarily those of the Federal Reserve Bank of Atlanta or the Federal Reserve System. Any remaining errors are the authors' responsibility.

Please address questions regarding content to Scott Baier, Assistant Professor of Economics, John E. Walker Department of Economics, Clemson University, Clemson, South Carolina 29634-1309, (864) 656-4534, and the Federal Reserve Bank of Atlanta, [sbaier@clemson.edu](mailto:sbaier@clemson.edu); Sean Mulholland, Assistant Professor of Economics, Auburn University at Montgomery, Montgomery, Alabama 36124-4023, (334) 244-3989, [smulholl@mail.aum.edu](mailto:smulholl@mail.aum.edu); Chad Turner, Visiting Assistant Professor of Economics, Williams College of Business, Xavier University, Cincinnati, Ohio 45207, (513) 745-3062, [turnerc1@xavier.edu](mailto:turnerc1@xavier.edu); or Robert Tamura, Associate Professor of Economics, John E. Walker Department of Economics, Clemson University, Clemson, South Carolina 29634-1309, (864) 656-1242, and the Federal Reserve Bank of Atlanta, [rtamura@clemson.edu](mailto:rtamura@clemson.edu).

Federal Reserve Bank of Atlanta working papers, including revised versions, are available on the Atlanta Fed's Web site at [www.frbatlanta.org](http://www.frbatlanta.org). Click "Publications" and then "Working Papers." Use the WebScriber Service (at [www.frbatlanta.org](http://www.frbatlanta.org)) to receive e-mail notifications about new papers.

**INCOME AND EDUCATION OF THE STATES OF THE UNITED STATES:  
1840-2000**

In order to understand the relationship between long-run economic growth and the role of inputs into the production process a long time series is needed. For the states of the United States of America, there exists data on output production, population, and enrollment that can be employed to enlighten us on the nexus between educational attainment and income per worker in each state. These data, however, have not been organized in a manner that lends itself easily to economic analysis. To this end, this paper makes three contributions: (1) it introduces original annual measures of years of schooling and average years of experience in the labor force for each of the states of the United States, generally from 1840 through 2000, (2) it constructs original real state per worker output estimates for 1850, 1860, 1870, 1890 and 1910, and combines them with existing data for 1840, 1880, 1900 and 1920 and 1929 through 2000, (3) it estimates the return to schooling and experience over this period. We provide a long term perspective on the return to human capital accumulation. Furthermore, it captures the educational choices made by individuals (aggregated to the state level) over much of the history of the United States. We use data from the decennial censuses of the United States, Richard Easterlin's work on state income, *Historical Statistics of the United States: Colonial Times to 1970* as well as information contained in annual *Statistical Abstracts of the United States* to produce these estimates.<sup>1</sup> These data, aggregated to the level of state education and income, show that investments in schooling are quite productive; that is, the estimated return to a year of schooling for the average individual in a state ranges from 8 percent to 12 percent. This range is robust to various time periods and various estimation methods. Although not necessarily producing similar results, we view this work as complementary to the work of Mulligan and Sala-i-Martin (1997, 2000).<sup>2</sup> By census region, we also document the long-term enrollment trends in primary, secondary, and tertiary schooling as well as the patterns of income growth across regions. We show both within region and across region convergence.

The remainder of the paper is organized as follows: The next section provides the accounting

---

<sup>1</sup>While we would like to go all the way back to the establishment of the United States as a nation 1776 (1788 as a Constitutional Republic), the data do not appear to be easily available to researchers prior to 1840. We envision that the data exist in some form at the state level, typically in the form of Reports of the State Superintendent of Schools, but we have not investigated these potential sources at this time.

<sup>2</sup>Mulligan and Sala-i-Martin (1997,2000) construct two different state level human capital measures for the census years 1940-1990, inclusive. Our years of schooling human capital measure is highly correlated with theirs, averaging approximately 0.8. See Appendix D for more detail.

framework for calculating average years of schooling by state. We present in graphical and tabular form the results of these calculations by census region. Section III presents our measures of state output per worker. Similar to the results from our years of schooling calculations, we find that among the nine census regions, there have been systematic leaders and laggards. Section IV contains our estimates for the returns to schooling and the returns to potential job experience. We find that OLS estimates are quite robust to alternative specifications, and that a year of schooling returns about 10 percent to an individual in additional productivity. Section V concludes and describes future work.

## II. EDUCATION IN THE STATES

In this section we present average schooling measures for each of the nine census regions.<sup>3</sup> We present our methodology for calculating years of schooling for the average labor force participant in each state.<sup>4</sup> We also compare each labor force regional average with the labor force average for the US. Rather than presenting graphs with 50 lines or tables with 50 rows, aggregation at the census region is a parsimonious manner to present the data.<sup>5</sup> Later in the empirical sections, we use the data for each state.

We use a perpetual inventory method, employed by Barro and Lee (1993) and Baier, Dwyer and Tamura (2004) for cross country tabulations, in order to construct average years of schooling in the labor force for each state. Because we are interested in output per worker, it is more appropriate to calculate the average years of schooling in the labor force instead of the average years of schooling of all state residents.<sup>6</sup> We also are unable to account for changes in the labor force participation rates by educational category, because we do not have any historical data on labor force participation by education category prior to 1960.

---

<sup>3</sup>For a listing of states within each region, see Appendix A.

<sup>4</sup>Additional details on the derivation and the data sources are furnished in Appendix B.

<sup>5</sup>We do present information about maximum gaps between states in some of our tables.

<sup>6</sup>Ideally we would use information to produce average years of schooling for men and women separately in the labor force, however, enrollment information by sex is not consistently available. However Series H 433-441, page 370 of *Historical Statistics of the United States: Colonial Times to 1970*, indicates that there was little difference

	sex	1850	1860	1870	1880	
in enrollment rates of men and women:	male	49.6	52.6	49.8	59.2	. From 1890 onward differences in
	female	44.8	48.5	46.9	56.5	

enrollment rates were less than one percentage point. We acknowledge that our calculations implicitly assumes that the labor force participation rate is common across men and women.

We assume that there are four categories of workers, those with no schooling (none), those exposed to primary schooling and no more (primary), those exposed to secondary schooling and no more (secondary), and those with exposure to higher education (college). Our enrollment data includes both public and private primary schools, secondary schools and institutions of higher education.<sup>7</sup> To calculate our average years of schooling, we assign the average years of schooling attained for each of these categories, with the uneducated group getting zero years of schooling. Suppressing the state subscript,  $H_t^i$  is the number of workers in the labor force in year t in education category i. The perpetual inventory method produces the following law of motion of these variables.

$$H_{t+1}^i = H_t^i (1 - \delta_t^i) + I_t^i, \quad i = \text{none, primary, secondary, college} \quad (1)$$

where  $\delta_t^i$  is the departure rate from the labor force between year t and t+1 and  $I_t^i$  is the gross flow of new workers into the labor force from education category i.

We assume three different departure rates: one for college workers,  $\delta_t^{\text{college}}$ , one for secondary workers,  $\delta^{\text{secondary}}$ , and one for all other workers,  $\delta_t^{\text{primary}}$ .<sup>8</sup> We assume these different rates for two reasons: (1) because a common rate produces a 2000 share of workers with some college significantly below the 50 percent reported in the census and (2) when we use a common departure rate for secondary, primary, and no education workers, we observe states where the fraction of the labor force exposed to elementary schooling is less than zero.

Although values of  $\delta_t^i$  are not directly available, we are able to calculate the departure rates using the following three part solution. First, we assume that workers with some college exposure do not disappear at a calculated rate, but only after 45 years of employment. Thus for college exposed workers, the law of motion becomes:

$$H_{t+1}^{\text{college}} = H_t^{\text{college}} - I_{t-45}^{\text{college}} + I_t^{\text{college}} \quad (2)$$

Dividing through by labor force in period t+1 and defining  $h_t^i$  to be the share of the labor force in year t in education category i produces:

$$h_{t+1}^{\text{college}} = h_t^{\text{college}} \frac{L_t}{L_{t+1}} - \frac{I_{t-45}^{\text{college}}}{L_{t+1}} + \frac{I_t^{\text{college}}}{L_{t+1}} \quad (3)$$

For the very early years,  $I_{t-45}^{\text{college}}$  is approximated using the first observed measure of higher education

---

<sup>7</sup>See Appendix B for details on the information.

<sup>8</sup>We deliberately omit the time subscript on the departure rate for the secondary education category. Our reasoning is discussed in greater detail later in this section. Also, we use a common departure rate for the primary and none educational categories, which we denote  $\delta_t^{\text{primary}}$ .

enrollment rates in  $t$ .<sup>9</sup> Finally we assume that college aged individuals are those between the ages of 18 and 24, inclusive. We assume that population in this age category is uniformly distributed between these ages, and that higher education enrollment rates are constant over these ages. Thus we assume that:  $I_t^{\text{college}} = r_t^{\text{college}} lfp_t^{\text{college}} \ell[18-24]/7$ , where  $r_t^{\text{college}}$  is the higher education enrollment rate,  $lfp_t^{\text{college}}$  is the labor force participation rate of college exposed individuals, and  $\ell[18-24]$  is the population of 18 to 24 year olds, inclusive. Once enough years have past, we use our own calculations for  $I_{t-45}^{\text{college}}$ .

The second part of our solution is to determine a departure rate for workers exposed to secondary schooling. Initially, we included the secondary education exposed workers in a category along with those workers exposed to elementary education and no education. However, we find that this results in calculated shares exposed to elementary education that are less than zero. As a result, we choose  $\delta^{\text{secondary}}$  for each state by matching the calculated shares of workers exposed to secondary education to those observed in the census years from 1940-2000. We note that unlike the departure rates for other educational categories,  $\delta^{\text{secondary}}$  is time invariant. For values, see Appendix B.

Although we are unable to calculate the departure rate for the remaining educational classes directly, the final step allows us to isolate the departure rate for the remaining educational classes,  $\delta_t^{\text{primary}}$ , using the following identity on the law of motion of the labor force:

$$L_{t+1} = H_{t+1}^{\text{college}} + H_{t+1}^{\text{secondary}} + H_{t+1}^{\text{primary}} + H_{t+1}^{\text{none}} \quad (4)$$

Using (1) and (2) to substitute out for each education category produces:

$$\begin{aligned} L_{t+1} = & H_t^{\text{college}} - I_{t-45}^{\text{college}} + I_t^{\text{college}} + H_t^{\text{secondary}} (1 - \delta^{\text{secondary}}) + I_t^{\text{secondary}} \\ & + H_t^{\text{primary}} (1 - \delta_t^{\text{primary}}) + I_t^{\text{primary}} + H_t^{\text{none}} (1 - \delta_t^{\text{primary}}) + I_t^{\text{none}} \end{aligned} \quad (5)$$

Dividing through by  $L_{t+1}$ :

$$\begin{aligned} 1 - h_{t+1}^{\text{college}} = & h_t^{\text{secondary}} \frac{L_t}{L_{t+1}} (1 - \delta^{\text{secondary}}) + (h_t^{\text{primary}} + h_t^{\text{none}}) \frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}}) \\ & + \frac{I_t^{\text{secondary}} + I_t^{\text{primary}} + I_t^{\text{none}}}{L_{t+1}} \end{aligned} \quad (6)$$

In order to get estimates of the flows into each education category, we use the following information:

---

<sup>9</sup>This is not much of an issue in the early years because higher education enrollments are near zero. Further details are discussed in Appendix B.

$$I_t^{\text{college}} = \frac{r_t^{\text{college}} lfpr^{\text{college}} \ell[18-24]_t \Theta}{7} \quad (7)$$

$$I_t^{\text{secondary}} = \frac{(r_t^{\text{secondary}} - r_t^{\text{college}} \Theta) lfpr^{\text{secondary}} \ell[14-17]_t}{4} \quad (8)$$

$$I_t^{\text{primary}} = \frac{(r_t^{\text{primary}} - r_t^{\text{secondary}}) lfpr^{\text{primary}} \ell[5-13]_t}{9} \quad (9)$$

$$I_t^{\text{none}} = \frac{(1 - r_t^{\text{primary}}) lfpr^{\text{primary}} \ell[5-13]_t}{9} \quad (10)$$

where as before in year  $t$   $r_t^i$  is the enrollment rate in education category  $i$ ,  $lfpr_t^i$  is the labor force participation rates for each educational category, and  $\ell[i-j]_t$  is the population in age category  $[i-j]$ , inclusive.<sup>10</sup> The constant  $\Theta$  is an adjustment for the fact that, unlike primary and secondary schooling, there is no schooling level above the higher educational category; freshman college enrollment rates are much higher than sophomore enrollment rates.<sup>11</sup> Notice that we maintain the assumption of a uniform age distribution within age category and uniform enrollment rates within an age category. Combining (7)-(10) with (6) produces our estimate of the  $\frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}})$  term:

$$\frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}}) = \frac{1 - h_{t+1}^{\text{college}} - h_{t+1}^{\text{secondary}} - \left( \frac{I_t^{\text{primary}} + I_t^{\text{none}}}{L_{t+1}} \right)}{(h_t^{\text{primary}} + h_t^{\text{none}})} \quad (11)$$

Thus for the share of labor force with primary schooling exposure we produce:

$$h_{t+1}^{\text{primary}} = h_t^{\text{primary}} \frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}}) + \frac{I_t^{\text{primary}}}{L_{t+1}}, \quad (12)$$

---

<sup>10</sup>For labor force participation rates we used data from the 1940-2000 censuses to determine average labor force participation rates by educational attainment. We use .91, .82 and .60 for  $lfpr^{\text{college}}$ ,  $lfpr^{\text{secondary}}$ , and  $lfpr^i$ ,  $i = \text{primary, none}$ . We used these labor force participation rates for the entire 1840-2000 period. While it may seem strange to use a constant labor force participation rate, in 1840 the labor force participation rate for 14-65 year old individuals was 51 percent and in 1900 the labor force participation rate for this same category was 57 percent. Since the majority of our labor force is either without education or with only primary education in this period, we feel that holding labor force participation rates constant over time across education categories is reasonable.

<sup>11</sup>The fact that the conditional probability of attending *increases* after the second year of higher education with years attended exacerbates this problem. We chose  $\Theta$  in order to best fit both the higher education share as well as the secondary schooling share for each state. See appendix B for the values of  $\Theta$  for each state.

and we then use the following adding up restriction for the none share:

$$h_{t+1}^{\text{none}} = 1 - h_{t+1}^{\text{college}} - h_{t+1}^{\text{secondary}} - h_{t+1}^{\text{primary}}.^{12} \quad (13)$$

We use information from the 1940-2000 Censuses to get estimates for expected number of years of schooling completed, conditional on being in each education category for each state. These expected years of schooling by category are represented by  $yr s_{it}^{\text{college}}$ ,  $yr s_{it}^{\text{secondary}}$ , and  $yr s_{it}^{\text{primary}}$ . For the intervening years we log linearly interpolate. Initial values for  $yr s_{it}^{\text{college}}$ ,  $yr s_{it}^{\text{secondary}}$ , and  $yr s_{it}^{\text{primary}}$  are set at 4, 10 and 14 for primary, secondary and higher education, respectively, in the year that data becomes available for each state.<sup>13</sup> We then log linearly interpolate from these initial values to the 1940 value. Thus for state  $i$  we calculate average years of schooling in the labor force as:

$$\widehat{E}_{it} = h_{it}^{\text{college}} yr s_{it}^{\text{college}} + h_{it}^{\text{secondary}} yr s_{it}^{\text{secondary}} + h_{it}^{\text{primary}} yr s_{it}^{\text{primary}} \quad (14)$$

To account for interstate migration, we adjust our years of schooling measure by residents state of birth reported in the 1850 through 2000 Censuses.<sup>14</sup> We assume that all education is undertaken in an individual's state of birth and that all current migrants are educationally representative of their birth state. Due to data limitations, our assumptions do not allow for selective migration. Let  $\widehat{E}_{jt}$  be the years of schooling at time  $t$  for those born in state  $j$ . Our estimate of years of schooling in state  $i$  therefore is:

$$E_{it} = \sum_{j=1}^{52} S_{ijt} \widehat{E}_{jt} \quad (15)$$

where  $S_{ijt}$  is the share of state  $i$  residents in year  $t$  that were born and educated in state  $j$ . There are 52 categories: 50 states, the District of Columbia, and the foreign born. For foreign born we assume that the individuals come from the  $k^{\text{th}}$  percentile of the primary, secondary and higher education distributions. We use the information from each of the 1940-2000 Censuses to determine the best

---

<sup>12</sup>There are occasions when  $h_t^{\text{none}} < 0$ . In these instances, we set  $h_t^{\text{none}} = 0$  and renormalize the shares to sum to 1. These instances are rare and small in absolute value.

<sup>13</sup>See Appendix B for more details on the various values of average years of schooling.

<sup>14</sup>In 2000, data availability is limited. The census reports the fraction of a state's residents that were born in that state,  $S_{ii}$ , and the fraction that is foreign born  $S_{i,for}$ . However, for those residents of a state who were not born in that state ( $S_{ij}$ ,  $j \neq i$ ,  $j \neq for$ ), only the census region of birth is given. Conditioned on living in state  $i$  and being born in census region  $k$ , we assume the probability of having been born in state  $j$  is equal the population of state  $j$  divided by the population of region  $k$ . We make the necessary adjustment when the region of birth contains the state of residence. As data is not available for 1840, we assume the shares in 1840 are identical to the values in 1850. Also, data is not available for Alaska and Hawaii in 1940 and 1950. We assume these shares are identical to the values in 1960.

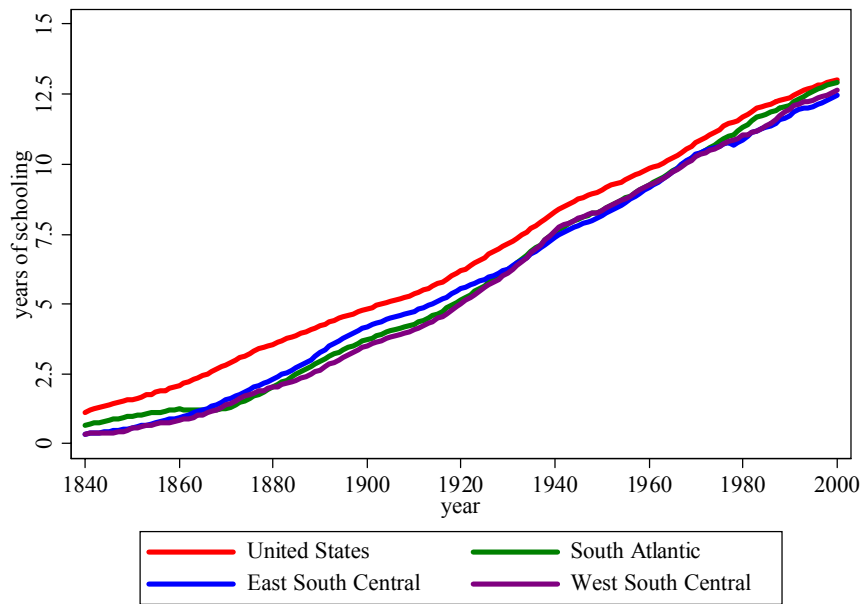
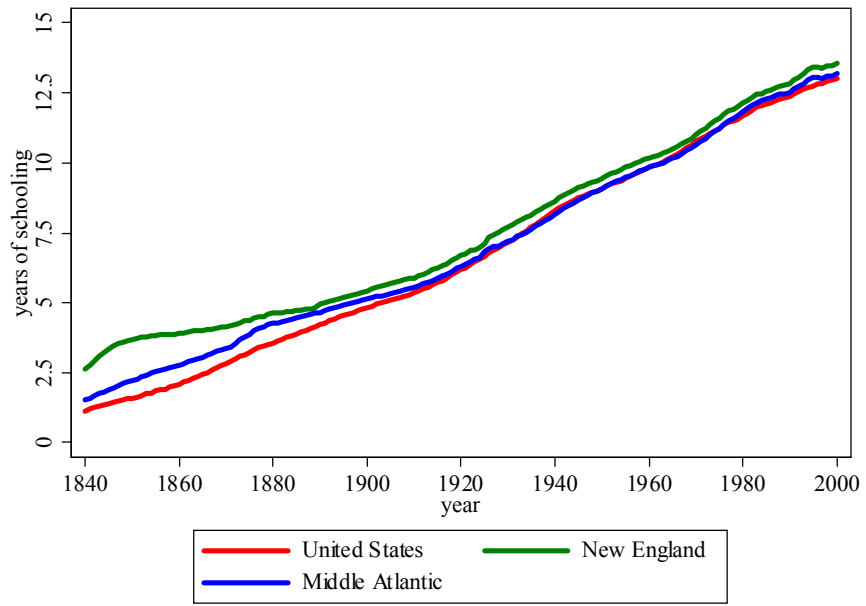


fitting  $k^{th}$  percentile for each state and census year in order to match the state's average years of schooling. For the non census years between 1940-2000 we linearly interpolate the shares born in state  $j$  residing in state  $i$  in year  $t$ . For years prior to 1940 we assume that foreign born workers have the average  $\bar{k}^{th}$  percentile, where the average is for the 1940-2000 period, and is state specific.<sup>15</sup>

To illustrate our years of schooling measure, the next four figures display the average years of schooling in the labor force by census region. While initial conditions certainly come into play in the first few years, within 20 years, the initial conditions have little impact. Thus New England, the Middle Atlantic and Pacific regions were clearly education leaders in the US. All three regions remain above the average years of schooling in the US throughout the entire 1840 to 2000 period. Figure 4 indicates that the East North Central and, by 1880, the West North Central were educational leaders as well. From 1880 to 2000 the labor forces of these five regions were better educated than the average person in the labor force in the US. In contrast, the South Atlantic, East South Central and West South Central regions were educational laggards. They start with less schooling than the average in the US, and remain below average throughout the data. However by 2000, these three regions have closed the gap between themselves and the US. Figure 3 illustrates the different behavior of the Mountain region. While the Pacific region remained above the US average, the Mountain region initially lagged behind the US, and in fact lagged behind the southern states from roughly 1850 to 1870. However from 1920 to the present the Mountain region was either at or above the US average in schooling. These results are summarized in Table 1 below.

---

<sup>15</sup>Details are in Appendix B. For information on how well our measure matches the Census data from 1940 to 2000 see Appendix D.



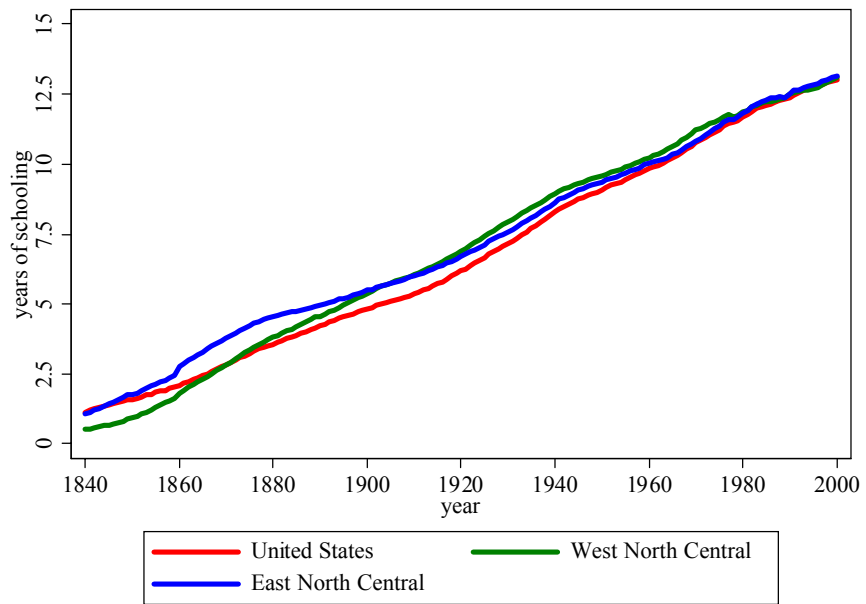
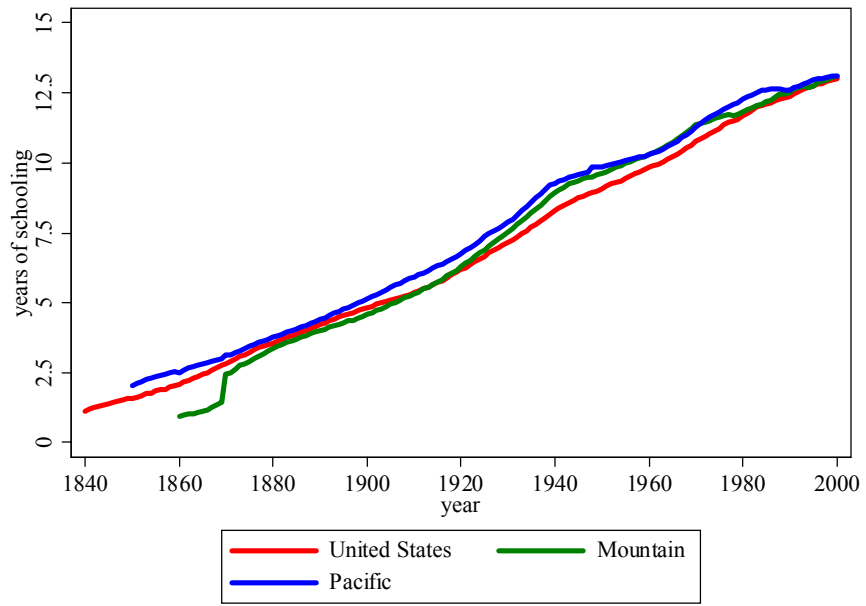


Table 1: Average Years of Schooling in the Labor Force

	1840	1860	1880	1900	1920	1940	1960	1980	2000
United States	1.14	2.09	3.56	4.83	6.18	8.28	9.83	11.7	13.0
New England	<b>2.61</b>	<b>3.90</b>	<b>4.62</b>	5.42	6.69	8.64	10.2	12.1	<b>13.5</b>
Middle Atlantic	1.52	2.76	4.25	5.15	6.30	8.15	9.84	11.8	13.2
South Atlantic	0.66	1.24	2.03	3.74	5.14	7.57	9.28	11.3	12.9
E. South Central	0.36	0.95	2.31	4.20	5.56	7.38	9.17	10.9	12.4
W. South Central	0.36	0.84	2.01	3.52	5.02	7.61	9.24	11.0	12.6
Mountain	-	0.95	3.37	4.58	6.29	8.94	<b>10.3</b>	11.8	13.1
Pacific	-	2.52	3.74	5.13	6.75	<b>9.25</b>	10.3	<b>12.3</b>	13.1
W. North Central	0.52	1.82	3.80	5.37	<b>6.88</b>	8.93	10.2	11.9	13.1
E. North Central	1.08	2.76	4.54	<b>5.48</b>	6.69	8.62	10.0	11.8	13.1
max. region gap	2.25	3.07	2.60	1.97	1.86	1.87	1.13	1.41	1.10
state max.	3.10	4.60	5.26	6.10	7.41	10.3	11.2	12.5	14.1
state min.	0.24	0.51	1.08	2.62	3.78	6.25	8.28	10.7	12.2

Table 1 contains the labor force weighted average years of schooling for each of the nine census regions and the average for the US for various years. For the US as a whole, the typical worker in 1940 had completed primary schooling and a quarter year of high school. By 1980 the typical worker was just about a high school graduate. In 2000 the labor forces in all regions have average schooling above 12 years. In 1880 the maximum gap between regions, 2.6 years, existed between the New England and West South Central regions. We pick 1880 as this is likely to be the first year in which initial conditions have no effect on the estimates. By 1900 the maximum gap between regions dropped to 1.97 years and existed between the East North Central and West South Central regions. From 1900 to 2000 the educational gap continues to narrow, reaching a nadir of 1.10 years in 2000.

Table 2 presents the maximum gap between regions, in the row marked R, and states, in the row marked S, at the decadal frequency, since 1890. Table 2 illustrates the clear convergence across regions, except for the very end. The evidence for the states is also compelling. The third row of Table 2, marked  $\widehat{S}$  contains the maximum gap between the 50 states of the US, dropping the District of Columbia.<sup>16</sup> We suppress the information where there is no change in the results with

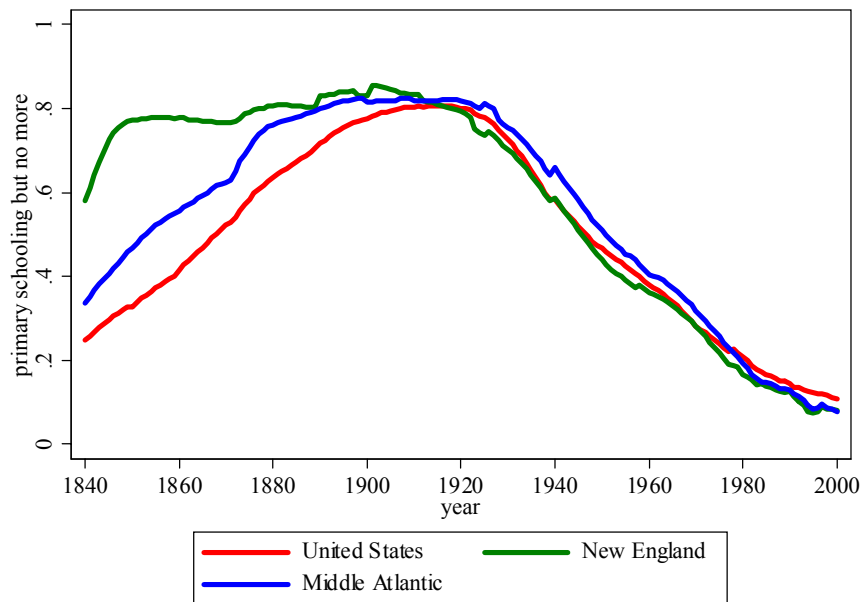
<sup>16</sup>The District of Columbia poses challenges in the latter years of the data. Specifically secondary and higher education enrollment rates in D.C. are exceedingly high. See Appendix B for how we dealt with “excess” enrollments

and without D.C. As can be seen, the same pattern arises with an increase in maximum gap in years of schooling from 1990 to 2000.

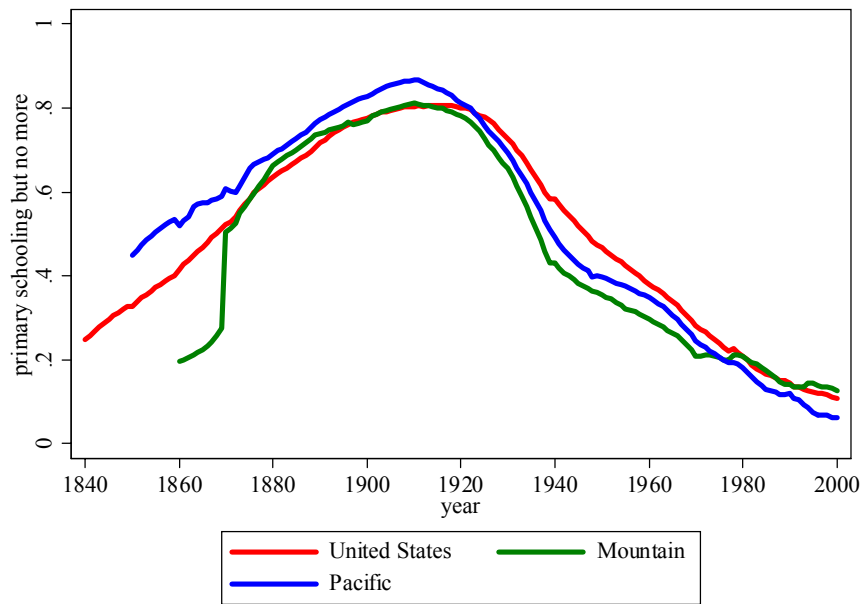
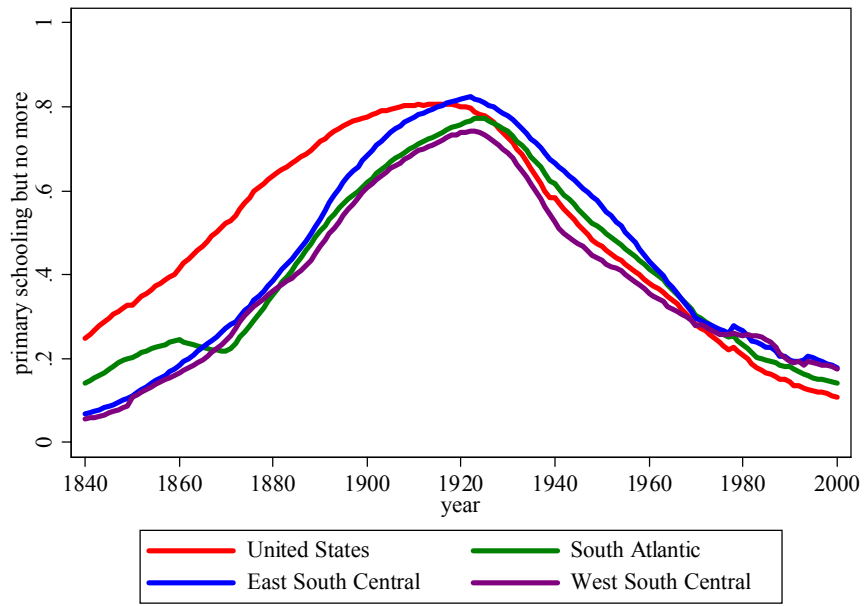
Table 2: Maximum Schooling Gaps between Regions and States

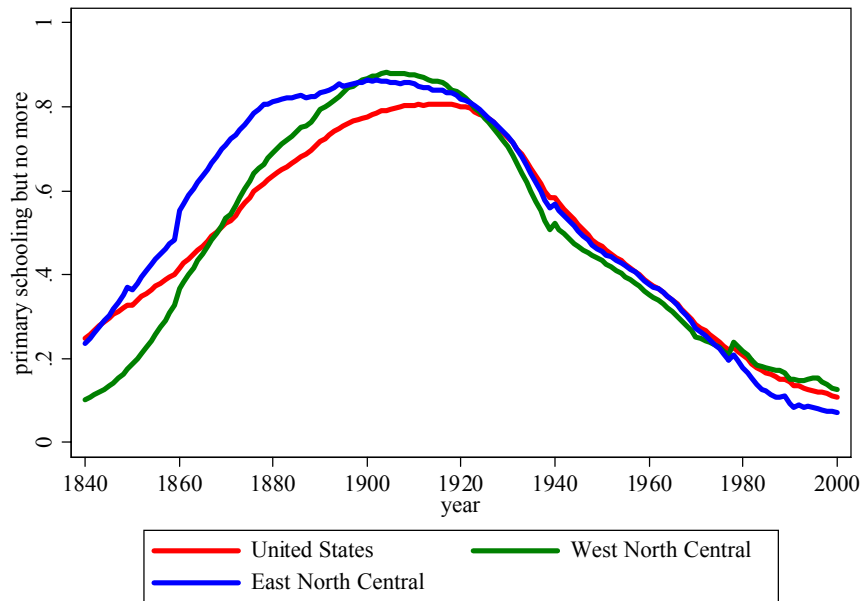
	1890	1900	1910	1920	1930	1940	1950	1960	1970	1980	1990	2000
R	2.30	1.97	1.98	1.86	1.84	1.87	1.73	1.13	1.07	1.41	1.05	1.10
S	3.79	3.48	3.69	3.63	3.88	4.06	3.34	2.91	2.21	1.79	1.79	1.95
$\hat{S}$										1.71	1.52	1.53

The differences in average years of schooling between regions are the result of systematic differences in enrollment rates across regions. New England, Middle Atlantic, Pacific, East North Central and, with a short lag, West North Central regions led the nation in educational attainment. These regions were the first to provide universal primary schooling, universal secondary schooling, and near universal higher education. In contrast, the South Atlantic, East South Central and West South Central regions lagged behind the country in each of these education categories. Finally the Mountain region is in between these two extreme groups. The next four figures illustrate the average fraction of the labor force that has been exposed to primary school, but no more.



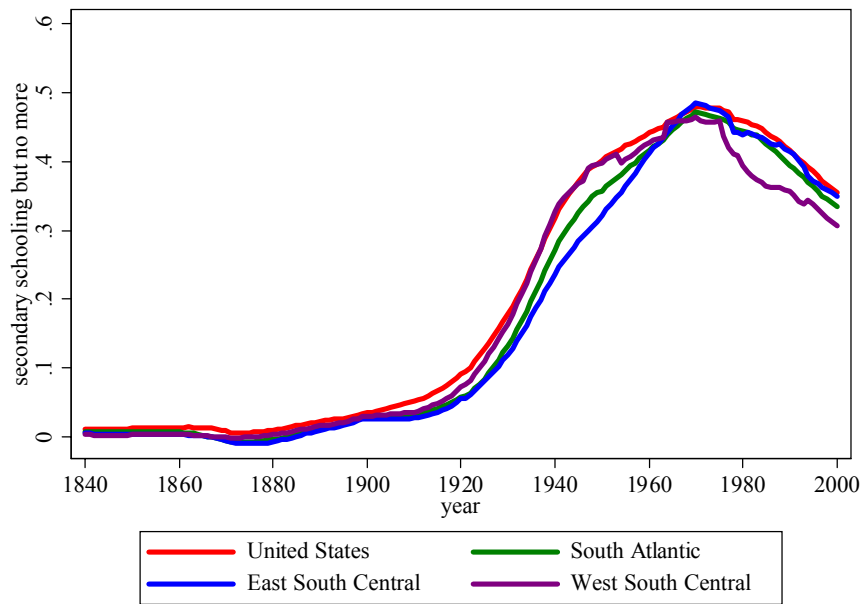
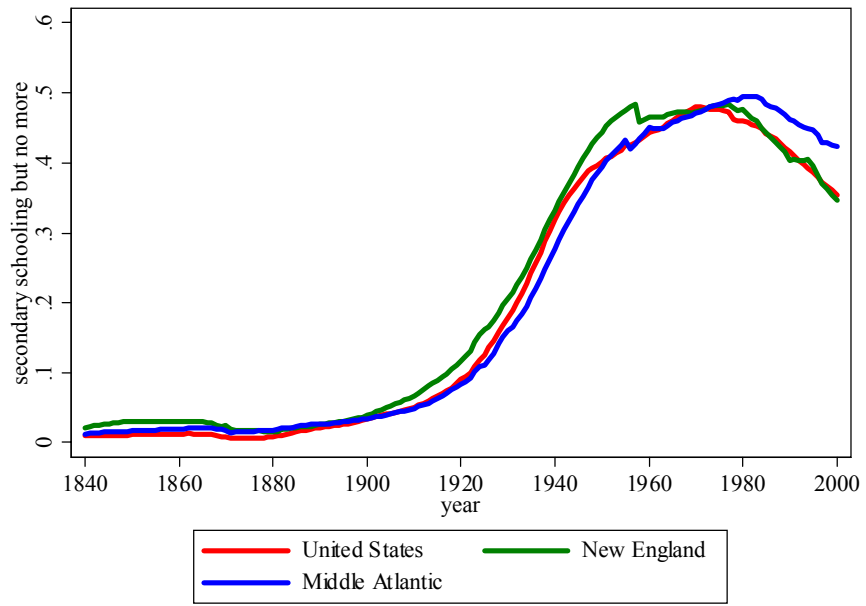
in D.C.



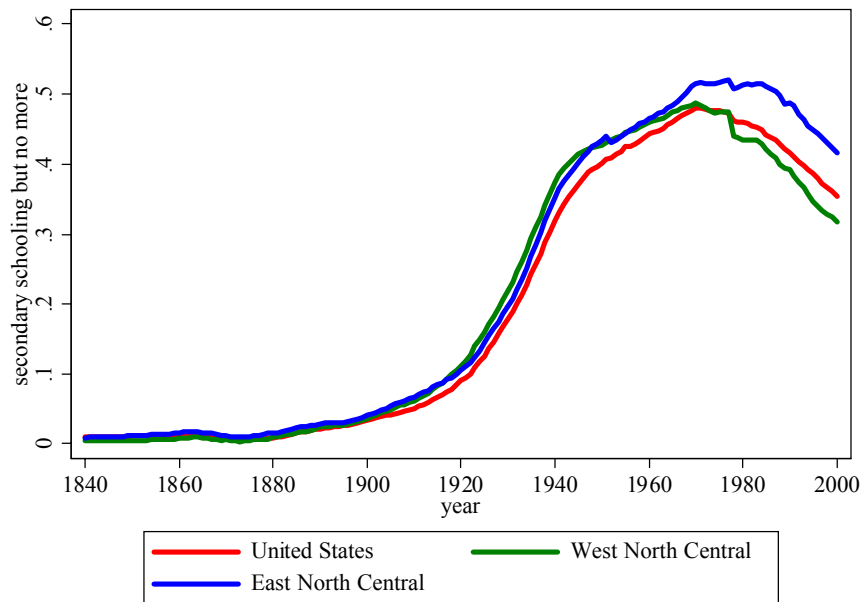
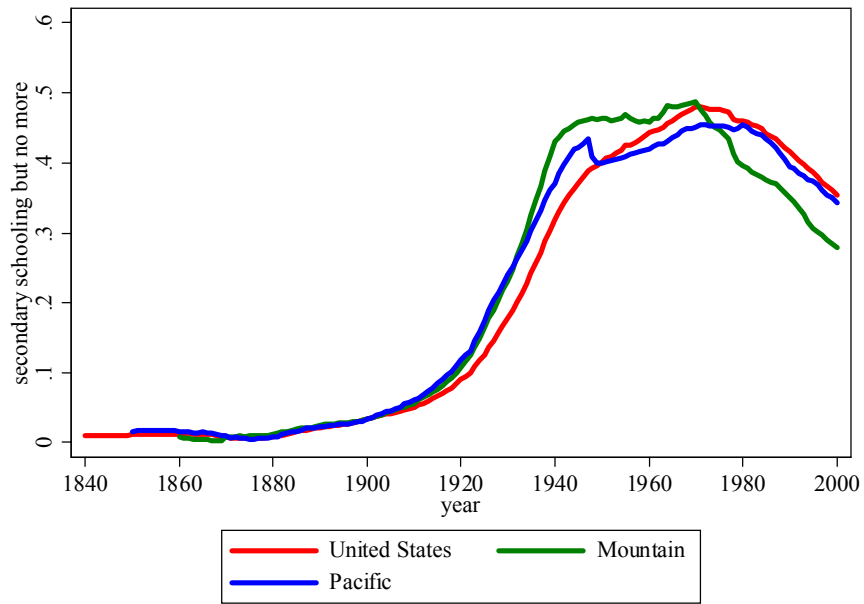


As the previous four figures illustrate, the South Atlantic, East South Central and West South Central regions display the lowest education exposure. From 1840 until about 1910 each of the three regions had a lower share of the labor force with elementary schooling exposure, and, as will be shown below, a lower fraction with secondary schooling exposure and higher education exposure as well. From 1920 to 2000 two of these regions have a greater share of the labor force with no more than an elementary schooling, and all three are higher after 1970. New England, Middle Atlantic, Pacific, East North Central and to a lesser degree the West North Central regions are educational leaders in the US. These regions led the nation in educating their residents first in primary school, then in secondary school and finally in higher education. The New England, Middle Atlantic, East North Central and to a slightly lesser degree the West North Central have higher share of the labor force with elementary schooling exposure than the national average from 1840 (roughly 1870 for the West North Central) until the early part of the 20th century, between 1900 and 1920.

The next four figures illustrate the evidence of some secondary schooling exposure but no more.



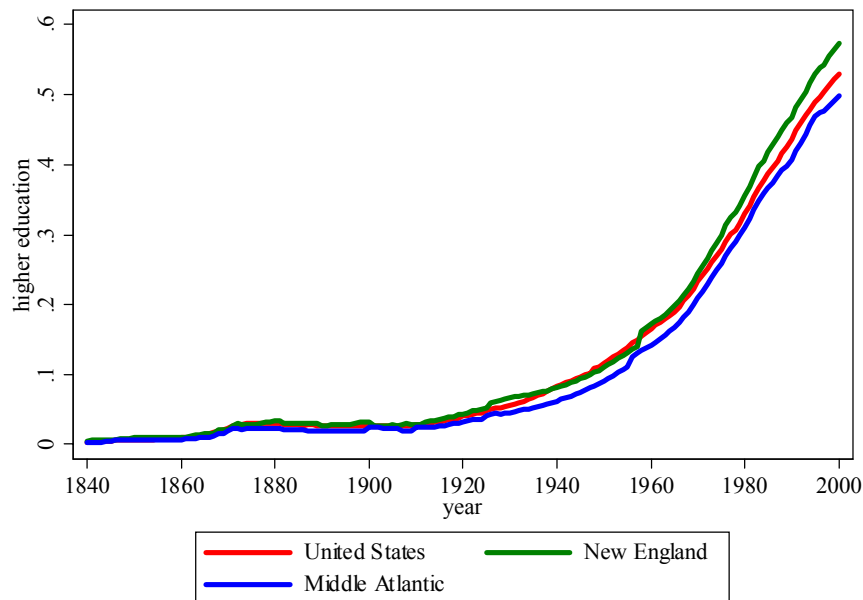


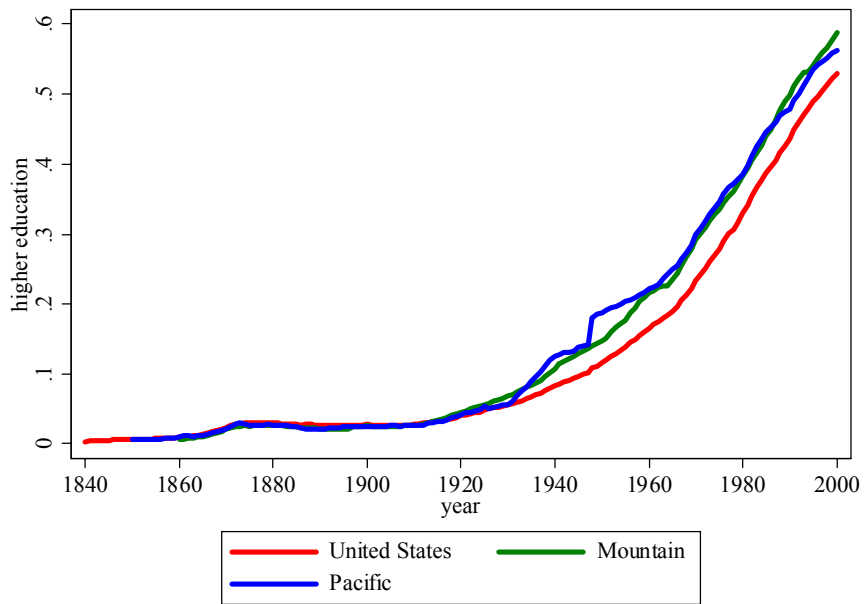
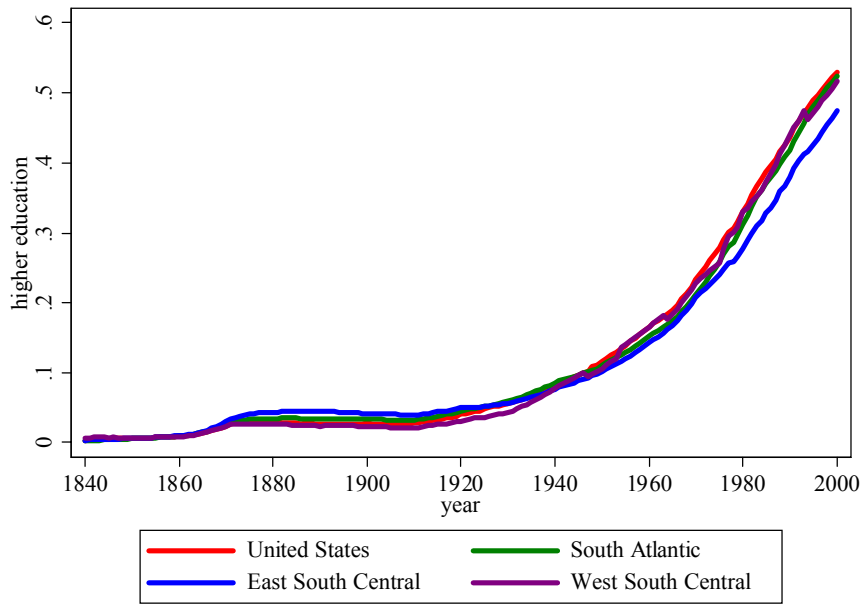


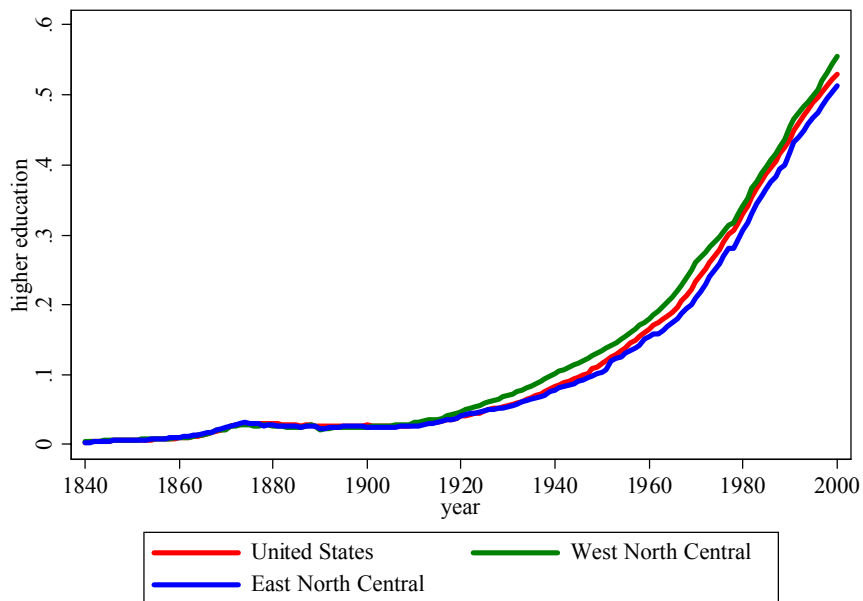
For secondary schooling exposure and no more, the nine census regions behave much like they did in elementary schooling exposure. From 1840 to 1940 the Pacific, and for 1840 to 1960, the Middle Atlantic, East North Central and West North Central regions display higher than average shares exposed to secondary schooling. As Goldin (1999) and Goldin and Katz (2000) have shown,

these were the leaders of the high school movement in the US as well as the world. The South Atlantic, East South Central, West South Central regions all lagged behind the average for the US from 1840 to this day. Combining these exposure rates with the primary exposure rates shows that the South Atlantic, East South Central, and West South Central clearly have the smallest portion of their labor force exposed to higher education.

The next four graphs present this the evidence for higher education. The regions with higher share of the labor force exposed to higher education are New England, West North Central, Mountain and Pacific. The South Atlantic, East South Central and West South Central regions remain below average throughout the entire time period. The Middle Atlantic and East North Central regions seem to almost mimic the national average.







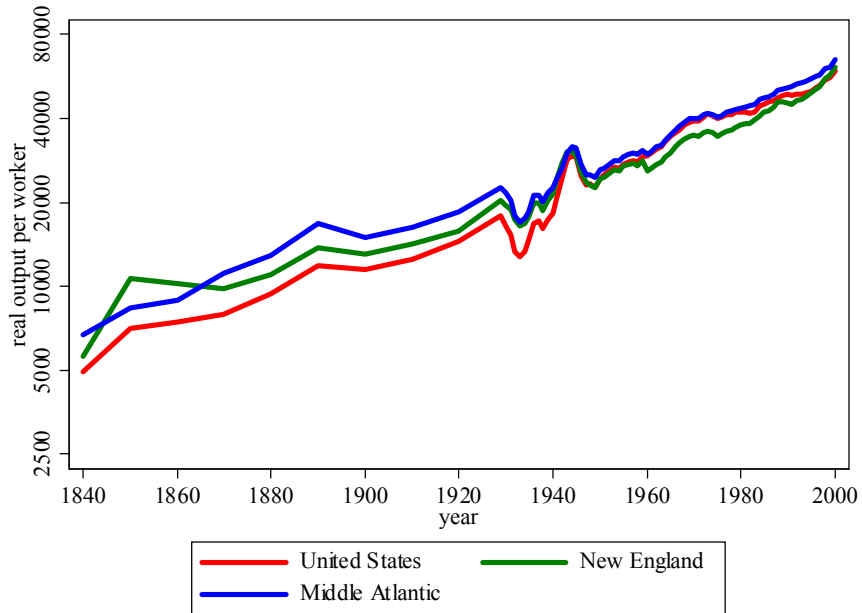
### III. STATE PER WORKER OUTPUT

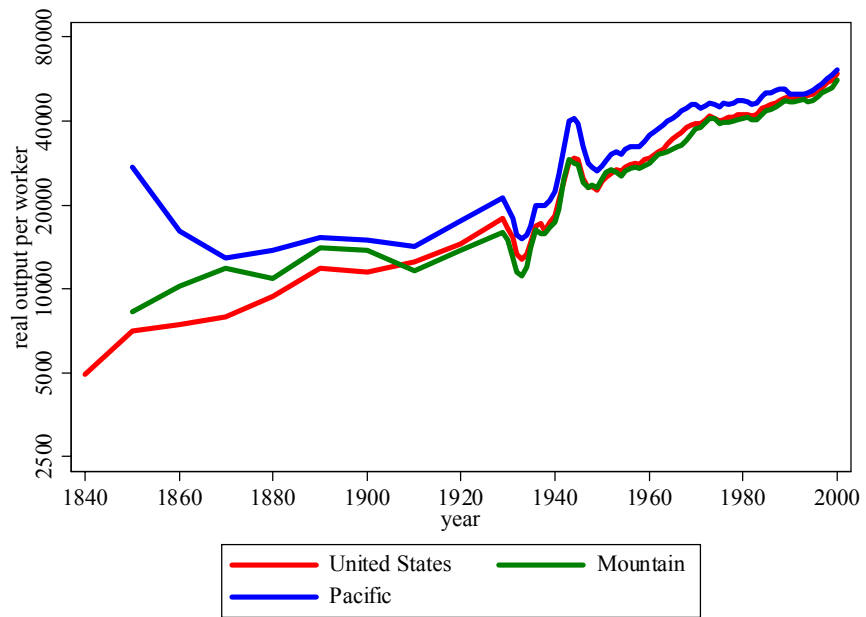
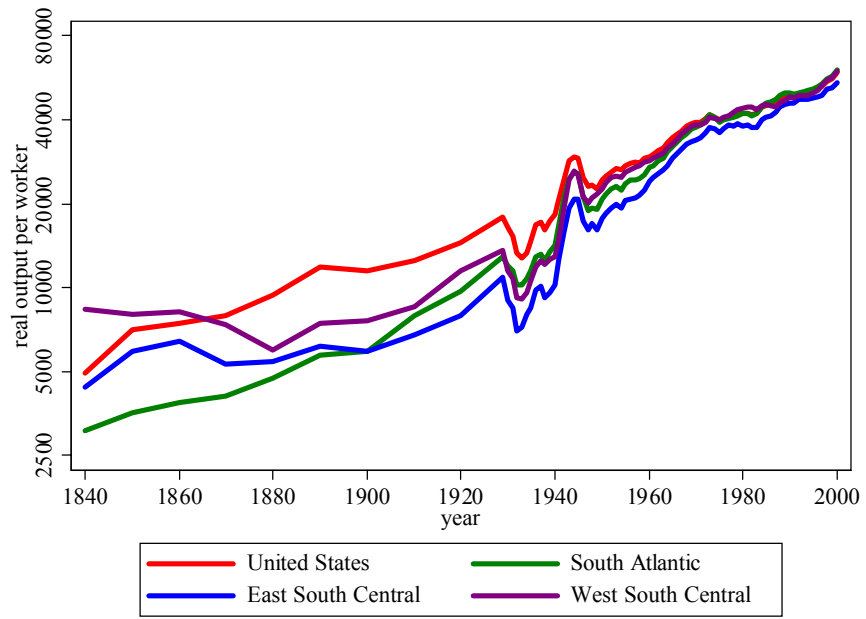
This section presents both original and existing data on state per worker output converted into real 2000 dollars.<sup>17</sup> In addition to the work of Easterlin (1960a,b), who provides per capita income in 1840, 1880, 1900, and 1919-1921 (1920), and government data from 1929-2000, we add our original estimates of state per capita income for 1850, 1860, 1870, 1890, and 1910. Our work uses government sources to produce estimates of real agricultural output, manufacturing output and mining output per state for these years. In combination with our measures of the labor force and the sectoral allocation of the labor force, we construct estimates of the non-agricultural, non-manufacturing non-mining output. With these estimates we create output per worker by state. The details of

<sup>17</sup>We convert all nominal values into real 2000 dollars, using the GDP deflator data from Gordon (1999) for years 1870-2000. For values between 1840-1869 we use the wholesale price index from the *Historical Statistics of the United States: Colonial Times to 1970* to compute inflation rates over this period. We then use the calculated wholesale price inflation to create a GDP deflator for the 1840-1869 period. To account for regional price differences, we use Berry, Fording, and Hanson (2000), Mitchener and McLean (1997), and Williamson and Linder (1980). The first deflators provide measures of output or income in constant national dollars and the regional price corrections adjust for regional price variation. For the 1840-1880 period we extrapolated the trend in relative price levels for the Mountain and Pacific region. Thus the output measures are best thought of as real income per worker. More details are available in Appendix B.

these calculations are in Appendix C. We note that the data from 1850-1920 are for state output per worker. For the period 1929-2000, the data are for state income per worker.

The next set of figures displays the regional average output per worker and the US average output per worker. As with the educational measures, we present the data in regional aggregates in order to easily facilitate data presentation. The real income per worker series has many similarities with the educational attainment data. The Middle Atlantic and Pacific regions are consistently more productive than the US from 1840-2000, and the South Atlantic, East South Central and West South Central regions are consistently less productive than the US from 1840-2000. The remaining three regions, Mountain, West North Central and East North Central are essentially as productive as the US from 1840-2000.





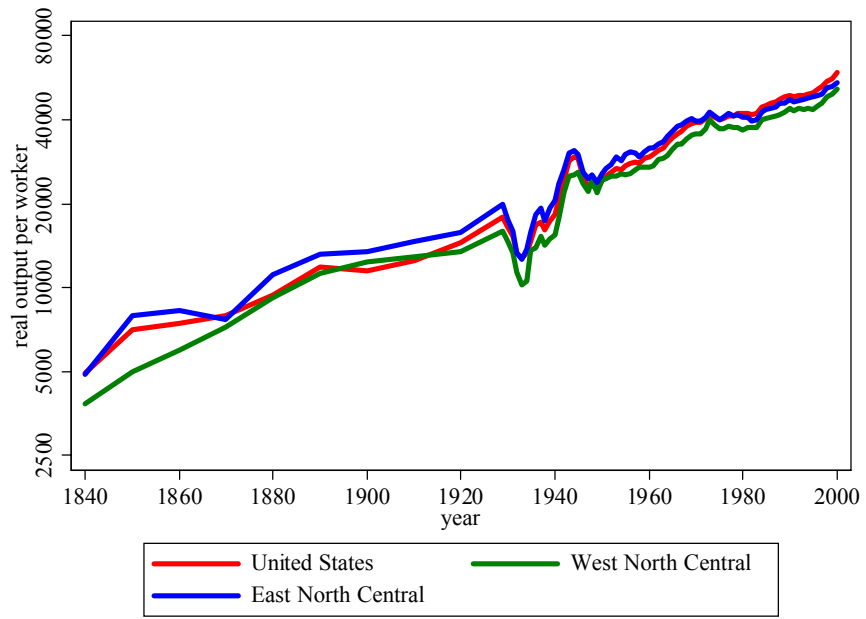


Table 3: Real Output per Worker

(regional leaders in bold)

	1840	1860	1880	1900	1920	1940	1960	1980	2000
United States	4950	7490	9448	11477	14430	18328	29514	42083	58791
New England	5640	10216	10998	13073	15706	21518	26042	38074	61426
Middle Atlantic	6709	8952	12954	14947	<b>18469</b>	<b>22639</b>	29854	43667	<b>64758</b>
South Atlantic	3089	3882	4751	5929	9770	14278	26982	42058	60216
E. S. Central	4391	6442	5447	5900	7947	10240	24092	37899	54134
W. S. Central	<b>8363</b>	8209	5971	7641	11512	12993	28521	43845	59833
Mountain	-	10236	10913	13838	13823	17247	28272	40690	56277
Pacific	-	<b>16167</b>	<b>13787</b>	<b>14992</b>	17607	22302	<b>35638</b>	<b>47185</b>	61374
W. N. Central	3825	5945	9248	12395	13497	15515	26991	36952	51527
E. N. Central	4867	8265	11147	13440	15841	20512	31641	40972	54162
region $\frac{\max}{\min}$	2.71	4.16	2.90	2.54	2.32	2.21	1.48	1.28	1.26
state max.	9218	16672	18972	17088	20492	28797	38531	62117	82438
state min.	2660	3144	3297	3678	6019	7135	20032	31558	41653
state $\frac{\max}{\min}$	3.47	5.30	5.75	4.65	3.40	4.04	1.92	1.97	1.98

As apparent in the figures as well as Table 3, real output per worker has increased substantially in the US, and across all regions. Consistent with evidence for the US from Baier, Dwyer and Tamura (2004), real output per worker grew at an annual rate of 1.6 percent per year. The nine census regions had annual real output per worker growth rates of 1.5 (New England), 1.4 (Middle Atlantic), 1.9 (South Atlantic), 1.6 (East South Central), 1.2 (West South Central), 1.2 (Mountain), 1.0 (Pacific), 1.6 (West North Central) and 1.5 (East North Central). The surprising values come from the West South Central, Mountain and Pacific. In the case of the West South Central, the high value in 1840 comes from Louisiana, with real output per worker of 9218 dollars. Workers in the only other state in this region for 1840, Arkansas, realized a real output per worker of 5313 dollars. From 1860 to 2000, the West South Central saw real output per worker grow at 1.4 percent per year. For the Mountain region, in 1860 only New Mexico and Utah are in the data. Each has worker productivity in excess of 9800 dollars compared with the US value of 7500 dollars. In 1870 Colorado, Montana and Nevada enter the data. Montana, New Mexico and Utah all have worker productivity of about 5900 dollars, however Colorado and Nevada are very productive mining states. These two states each have worker productivity in excess of 20,000 dollars. In 1880, Arizona and



Idaho arrive in the data; all but New Mexico and Utah have worker productivity in excess of the US average, 9447 dollars. In the case of the Pacific region, California, Oregon and Washington all have real output per worker values in excess of 10,000 dollars. These states were likely very high cost of living states as many manufactured goods would have to be imported from the rest of the US or abroad. Real output per worker for the Pacific region grows at an annual rate of 1.3 percent from 1880-2000 and 1.4 percent from 1900-2000. However over the last 80 years, 1920-2000, the Pacific region enjoyed real output per worker growth of 1.5 percent per year.

Our results are also consistent with those in Goldin and Margo (1992a). While they found falling real wages for artisans, laborers and clerical workers between 1840-1856, this is consistent with what we find for non agricultural workers. While agricultural workers saw rising output per worker from 1840 to 1860, 3984 dollars to 7937 dollars, their share of the labor force fell from 76 percent to 53 percent. Between 1840 and 1860 we find that nonagricultural workers real output falls from 8009 to 6986, or a decline of 0.7 percent per year.<sup>18</sup>

The effects of the Civil War are quite prominent in the figures, and are evident in Table 3. The states of the old Confederacy, South Atlantic, East South Central and West South Central clearly have lower growth rates. Between 1860 and 1880, these three regions experienced real annual income per worker growth of 1.0 percent, -0.8 percent and -1.6 percent, respectively. For the South Atlantic and East South Central, these understate the magnitude of the reduction in output per worker since 1870 is the nadir for these regions. The annual growth rates of income per worker from 1860 to 1870 for these three regions are 0.6 percent, -1.9 percent and -1.0 percent, respectively. In 1880 (1860) their relative worker productivity values were 50 (52) percent, 58 (86) percent and 63 (110) percent. By 2000 only the East South Central remains below the national average.

The final four rows of Table 3 present evidence on regional output per worker convergence. These contain the ratio of the maximum regional income per worker to minimum regional income per worker, the maximum and minimum state per worker income, and the ratio of the maximum state income per worker to minimum state income per worker. Inequality in 1870 and 1880 are certainly higher than in the pre Civil War period. Inequality in output per worker is reduced throughout the next century. By 1980 the relative region gap is about one third of its value in 1880, and the relative state gap is less than a third of its 1880 value. Though the relative state gap has barely increased

---

<sup>18</sup>Between 1840 and 1856 Goldin and Margo (1992a) present annualized real wage growth rates for artisans, laborers and clerical workers as: -0.7, 0.4 and 0. These average figures are obtained by equally weighting each of their four geographic regions.

somewhat in 2000 compared to its 1980 value, the relative region gap has fallen since 1880.<sup>19</sup>

#### IV. RETURNS TO SCHOOLING

Before we present evidence on the rate of return to schooling, it is necessary to deal with missing data on other inputs. Consider a model with two factors of production, human capital and all other inputs which we call physical capital. We assume production of a single final output is Cobb-Douglas. We assume perfect competition in factor markets and free mobility of capital. State  $i$  output per worker is given by:

$$y_{it} = A_{it} k_{it}^{\alpha} (\text{human capital})_{it}^{1-\alpha} \quad (16)$$

where  $k_{it}$  is physical capital per worker and  $\text{human capital}_{it}$  is human capital per worker. Under perfect competition in the output market, with final output as numeraire, the representative firm solves:

$$\max \left\{ A_{it} k_{it}^{\alpha} (\text{human capital})_{it}^{1-\alpha} - r_t k_{it} - w_t \text{human capital}_{it} \right\} \quad (17)$$

where  $r_t$  and  $w_t$  are the rental rate per unit of physical capital and human capital, respectively. Under competition firms choose physical capital in proportion to the human capital in the workforce:

$$k_{it} = \left( \frac{w_t}{r_t} \right) \left( \frac{\alpha}{1-\alpha} \right) \text{human capital}_{it} \quad (18)$$

Therefore substituting this back into the output equation produces:

$$y_{it} = A_{it} \left( \frac{w_t}{r_t} \left( \frac{\alpha}{1-\alpha} \right) \right)^{\alpha} \text{human capital}_{it} \quad (19)$$

We assume that  $\text{human capital}_{it}$  can be specified in a Mincerian fashion:

$$\text{human capital}_{it} = \exp(\beta E_{it} + \gamma x_{it}) \quad (20)$$

where  $E_{it}$  is years of schooling in state  $i$  in year  $t$ , and  $x_{it}$  is experience in state  $i$  in year  $t$ .<sup>20</sup> In order to construct average experience by state, we calculated average age in the state not enrolled

<sup>19</sup>These results are consistent with those found using state income per capita from 1880, 1900, 1920 and 1930-1990 at the decadal frequency in Tamura (2001).

<sup>20</sup>Those familiar with the standard Mincer earnings regression may wonder why we exclude the quadratic term in experience. This is because of aggregation bias, while one can construct a model in which the linear terms in education and experience are identified by state variation, the quadratic term is not identified upon aggregation. When we experimented with identification, the results confirmed the bias in estimation, and hence we ignore the diminishing returns to work experience. The results indicate that experience returns are significantly below that from additional schooling and hence suggest that ignoring the quadratic term is not problematic.

in school and under the age of 65. From average age we subtract the sum of our average years of schooling measure in the labor force and the 6 years before individuals traditionally begin school enrollment. With this definition of *human capital*<sub>it</sub> the “earnings regression” is:

$$\ln y_{it} = \ln A_{it} + \alpha \ln \left( \frac{w_t}{r_t} \left( \frac{\alpha}{1 - \alpha} \right) \right) + \beta E_{it} + \gamma x_{it} \quad (21)$$

However before we can estimate the return to a year of schooling, we must remember that we substituted the optimal physical capital per worker. Thus the actual return to schooling would be given by:

$$\text{return to schooling} = (1 - \alpha)\beta \quad (22)$$

Therefore we need an estimate of the share of output that labor receives,  $(1 - \alpha)$ . Table 4 provides evidence on the share of output received by human capital (labor) for early years.<sup>21</sup> Table 4 shows that human capital’s (labor’s) share of output is roughly between  $\frac{2}{3}$  and  $\frac{4}{5}$ . This seems to hold for very long periods of time.

---

<sup>21</sup>Lines (1)-(12) Table reprinted from Table 15, *National Income: A Summary of Findings*, Kuznets, NBER (1946), p. 50.

Lines (13)-(25) from Table 4, Denison, *The Sources of Economic Growth in the United States and the Alternatives Before US*, Committee for Economic Development (1962) p. 30.

Table 4: Labor Share and Capital Share of Income

line	period	<i>Emp.</i> <i>Comp.</i>	<i>Entrep.</i> <i>Net Inc.</i>	(1) + (2)	<i>Div.</i> (3)	<i>Int.</i> (4)	<i>Rent</i> (5)	(3) + (4) + (5)
		(1)	(2)					
1	1870-1880	50.0	26.4	76.5	15.8		7.8	23.6
2	1880-1890	52.5	23.0	75.4	16.5		8.2	24.6
3	1890-1900	50.4	27.3	77.7	14.7		7.7	22.4
4	1900-1910	47.1	28.8	75.8	15.9		8.3	24.2
5	1899-1908	59.5	23.8	83.3	5.3	5.1	6.4	16.7
6	1904-1913	59.6	23.3	82.9	5.7	5.1	6.3	17.1
7	1909-1918	59.7	23.3	83.0	6.5	4.9	5.7	17.0
8	1914-1923	63.0	20.8	83.8	5.6	5.3	5.3	16.2
9	1919-1928	65.1	18.3	83.4	5.4	6.0	5.2	16.6
10	1919-1928	61.7	19.5	81.2	5.6	6.1	7.1	18.8
11	1924-1933	63.1	16.6	79.7	6.5	7.8	5.9	20.3
12	1929-1938	64.9	15.9	80.8	6.6	8.4	4.3	19.2
13	1909-1913			69.5				30.5
14	1914-1918			67.0				33.0
15	1919-1923			69.5				30.5
16	1924-1928			69.7				30.3
17	1929-1933			69.2				30.8
18	1934-1938			70.4				29.6
19	1939-1943			72.1				27.9
20	1944-1948			74.9				25.1
21	1949-1953			74.5				25.5
22	1954-1958			77.3				22.7
23	1909-1958			71.4				28.6
24	1909-1929			68.9				31.1
25	1929-1958			73.0				27.0

We first estimate (22) on each year, because of the possibility of technological progress. If the returns to schooling and experience are constant over time, but Total Factor Productivity rises over time, i.e. rising  $A_{it}$ , then any inability to properly control for the rising level of TFP will induce an upward bias on our estimates to schooling. The following figure presents the annual variation of the returns to schooling with one standard error bands.<sup>22</sup> With only four exceptions the estimates are always positive, and with very few exceptions the estimates are at least two standard errors away from zero. It is clear from the figure that the rate of return to schooling fell dramatically during the Depression; from 1929-1936 our estimates do not differ statistically from 0. However from 1937-1959 rates of return to schooling exceed 8 percent and for 1942-1959 an additional year of schooling returns in excess of 11 percent in every year. These high rates of return correspond to the diminishing dispersion of education across states, as well as rising levels of schooling. Together these help to explain the Great Compression in the middle of the 20th century as identified by Goldin and Margo (1992b). The falling returns to an additional year of schooling from the 1950s through the 1970s is consistent with the work of Freeman (1976). Although muted due to aggregation, the rising returns to schooling in the latter half of the 1980s and the recovery from the 1990-1991 recession are consistent with those found in Murphy and Welch (1992).

---

<sup>22</sup>The results come from annual weighted regressions of log state output per worker on years of schooling and experience, where the weights are the labor force of each state. Over the entire period, 1840-2000, the mean return, inclusive of physical capital's return, to a year of schooling is .0949 with a mean standard error of .0385. From 1880-2000, the mean return, inclusive of physical capital's return, to a year of schooling is .1012 with a mean standard error of .0347.

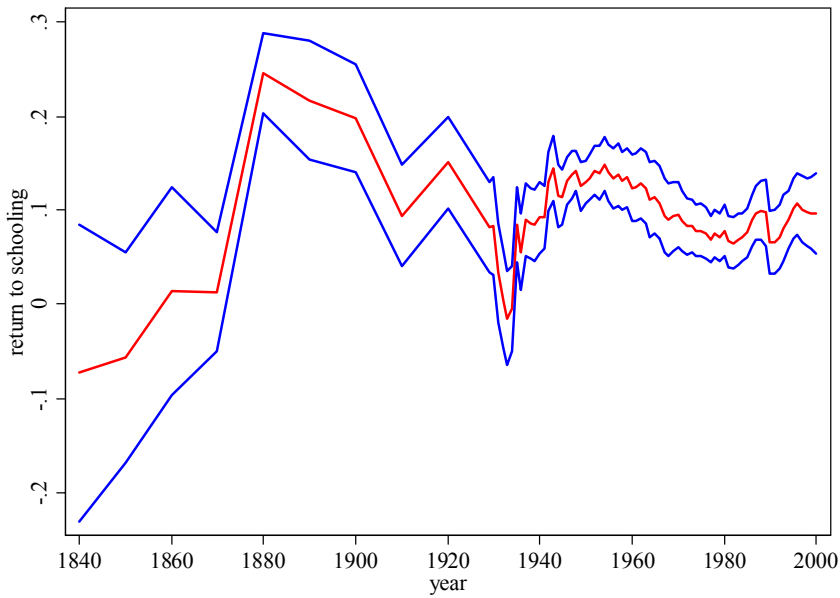


Table 5 contains the results of years of schooling regressed on real per worker output. The first three columns include year dummies to allow for more variation in technological change than a deterministic trend. The second column allows for a different return to schooling for Alaska. The third column allows for a different return to schooling and a different return to experience in Alaska. Under the hypothesis that TFP does not differ across states, i.e.,  $A_{it} = A_t$  for all  $i$ , differencing each state's log output per worker from the labor force weighted log US output per worker, years of schooling, and average experience from the labor force weighted US averages allows for the estimation of (Eq. 21) without any time controls. These differenced regressions are reported in the final three columns of Table 5.

Table 5: Earnings Regressions: Annual Data (standard errors)

$E$	.1244	.1213	.1213	.1239	.1226	.1214
	(.0045)	(.0043)	(.0043)	(.0045)	(.0043)	(.0043)
exp.	.0327	.0448	.0452	.0323	.0458	.0458
	(.0020)	(.0020)	(.0020)	(.0020)	(.0020)	(.0020)
$N$	4000	4000	4000	4000	4000	4000
$\overline{R}^2$	.8969	.9055	.9056	.2434	.2985	.3126
range	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]
year dummies	yes	yes	yes	no	no	no
ak $E$	no	yes	yes	no	yes	yes
ak exp.	no	no	yes	no	no	yes
differenced	no	no	no	yes	yes	yes

The results in Table 5 indicate an overall return to schooling, including the implied physical capital return, of 12 percent per year of schooling. These estimates imply that an additional year of schooling results in an 8 to 10 percent increase in worker productivity. These results are consistent with the evidence presented in Angrist and Krueger (1991), Staiger and Stock (1997), and Card (1995). The returns to experience, reflecting on-the-job training or learning by doing, are similar across all six columns. A one year increase in average experience raises worker productivity by three to four percent. At the individual level an additional year of experience returns between 2 percent to 4 percent in additional productivity.

Failing to account for the rising female labor force participation rate present over this period may result in poor estimates. To control for this we correct for the share of the labor force that is female (male) and interact these shares with average years of experience. This allowed us to separately measure the rate of return to experience for each sex. The results of these are contained in Table 6. The first three columns report the average return to schooling and average estimated returns to experience by sex with varying controls for Alaska. The remaining three columns are the differenced regressions, as in Table 5. The rows marked  $F$  and  $\text{Prob} > F$  contain the  $F$  statistic on the test of equality of returns to experience between men and women, and the p value of the statistic.

Table 6: Earnings Regressions: Annual Data (standard errors)

$E$	.1214	.1199	.1203	.1208	.1210	.1198
	(.0047)	(.0045)	(.0047)	(.0047)	(.0045)	(.0045)
exp male	.0361	.0456	.0469	.0359	.0443	.0468
	(.0028)	(.0027)	(.0027)	(.0027)	(.0027)	(.0027)
exp female	.0253	.0397	.0421	.0249	.0379	.0407
	(.0032)	(.0032)	(.0032)	(.0032)	(.0032)	(.0032)
$N$	4000	4000	4000	4000	4000	4000
$\overline{R}^2$	.8968	.9051	.9064	.2430	.3065	.3118
range	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]	[.08, .10]
year dummies	yes	yes	yes	no	no	no
$F$	5.61	1.83	1.24	6.01	2.14	2.01
Prob > $F$	.0179	.1757	.2663	.0143	.1437	.1560
ak $E$	no	yes	yes	no	yes	yes
ak exp	no	no	yes	no	no	yes
differenced	no	no	no	yes	yes	yes

The results of Table 6 indicate that the estimated returns to schooling are robust to the possible differences in returns to experience between men and women. It is reasonable to state that an additional year of schooling in a randomly chosen state returns 12 percent, and net of returns to physical capital, the typical worker would see an additional 8 to 10 percent increase in productivity. Rates of returns to experience for men and women are similar. In four of the six regressions we fail to reject the null that they are identical, only when we do not control for differential Alaskan returns to schooling do we reject the null. The typical worker becomes about 2 percent to 4 percent more productive at the individual level per additional year of experience. These results are almost identical to the estimates from Table 5.

One might be concerned that our estimates of the return to schooling may be biased because we assume a common intercept for all states in any time period. To address this concern, one way to correct for this is to allow for state specific effects. To help guide our think about alternative specifications that would correct for this potential bias, we return to equation (21)

$$\ln y_{it} = \ln A_{it} + \alpha \ln \left( \frac{w_t}{r_t} \frac{\alpha}{1 - \alpha} \right) + \beta E_{it} + \gamma x_{it} \quad (23)$$



The regression specification inferred from equation (21) is:

$$\ln y_{it} = c_i + b_t + \beta E_{it} + \gamma x_{it} + u_{it} \quad (24)$$

where  $c_i$  is the state specific fixed effects and  $b_t$  is a time specific effect common to all states. To correct for the state specific effects, there are two standard approaches to adjust for these effects: (1) standard fixed effects regressions or (2) the data can be first differenced. In both cases, it is required that there are no feedback effects from innovations in income influence future levels of educational attainment. If this is the case then standard fixed effects regressions or first-differencing will lead to inconsistent estimates of the return to schooling. To see if feedback effects are present, we follow Wooldridge (2002) and run a fixed effects regression with a lead of educational attainment in the specification. If the coefficient on educational attainment is statistically different from zero, then we will take this as evidence that contemporaneous innovations in income lead to future educational attainment.

Table 7: Fixed Effects with Leads of Education

	Earnings Reg
$E$	.0353 (.0310)
$E(t + 1)$	0.0692 (.0321)
exp	0.0428 (.0066)
$N$	663
Decade Dum.	yes

Table 7 reports the results from this specification. With time dummies we find that leads of education are correlated with contemporaneous income. As a result, the standard approaches to correct for state effects will lead to inconsistent estimates. To correct for the possibility of state specific effects, we difference the data in equation (21) to get

$$\Delta \ln y_{it} = \Delta b_t + \beta \Delta E_{it} + \gamma \Delta x_{it} + \Delta u_{it} \quad (25)$$

We are concerned that  $E(\Delta E_{it} \Delta u_{it}) = E[(E_{it} - E_{it-1})(u_{it} - u_{it-1})] \neq 0$  (which will be the case if there are feedback effects since  $E(E_{it} u_{it-1}) \neq 0$ ). To consistently estimate the above equation we

must find instruments for  $\Delta E_{it}$  that satisfy the standard instrumental variable assumptions; that is, (1) the instruments should be correlated with  $\Delta E_{it}$  and (2) the instruments are uncorrelated with the error term. Following Arraleno and Bond (1991) we could use lags of educational attainment in each state, but these lags failed all overidentification tests for all lags we attempted. Thus instead of using lags of educational attainment, we constructed a variable that may capture the changes in educational attainment related to regional convergence. More specifically, we create the variable

$$E_{it}^c = \left[ E_{it} - \frac{1}{N^R - 1} \sum_{j \neq i}^{N^R} E_{jt} \right] \quad (26)$$

where  $N^R$  is the number of states in region  $R$ . Thus, the variable  $E_{it}^c$  measures how far ahead or behind state  $i$  is relative to the rest of the states in the region. We use different lag levels and lagged growth rates of this variable to instrument for  $\Delta E_{it}$ . We report the results from three different lag structures, the first two satisfy the overidentification at the one percent significance level while the third produces similar point estimates but does not pass the overidentification test.<sup>23</sup>

Table 8: Earnings Regressions, Differenced Data IV Approach

	IV Educ	IV Educ	IV Educ
$E$	0.1486	0.1501	0.1413
	(.0229)	(.0223)	(.0222)
exp.	0.0161	0.0230	0.0274
	(.0153)	(.0149)	(.0148)
$N$	663	663	663
range	[.10, .12]	[.10, .12]	[.9, .11]
Decade Dum.	yes	yes	yes

Table 8 reports the results from these specifications. The returns to schooling are similar to the returns reported without attempting to control for state effects. The highest return is 12 percent and the lowest return is 9 percent. Thus, after controlling for state effects and allowing for income to influence educational attainment as in Bils and Klenow (2000), we still find the return to an additional year of schooling is roughly 10 percent.

<sup>23</sup>All other lag structures that produced similar results as far as the overidentification tests produced similar results for the return to schooling. The data indicate that lags of longer than three periods would fail the over identification tests and the returns to schooling would be much higher than 10 percent.

## V. CONCLUSION AND EXTENSIONS

This paper employs historical state enrollment and population data to produce original average years of schooling measures for each state from 1840 to 2000. We benchmark this measure to roughly match the census data in 1940 through 2000. We show that there has been tremendous increases in schooling in the US over the 1840-2000 period, with average years of schooling rising from 1 year to over 13 years. In addition there has been a reduction in the variance across states. We also construct original estimates for state per worker output for the census years 1850, 1860, 1870, 1890 and 1910. Coupling our constructed data with previous work by Easterlin and government data, we produce state per worker income measures for 1840 through 1920 at the decadal frequency and 1929 through 2000 at the annual frequency. We estimate rates of return to schooling for an individual and find that an additional year of schooling returned around 10 percent higher income. Our rate of return estimate is robust under alternative estimation methods.

Given this comprehensive measure of human capital accumulation, we envision future work combining these measures with historical data we have generated on state measures of physical capital. Using standard growth accounting methodologies, we will then estimate the contribution of aggregate input growth and total factor productivity growth on income growth across the United States. Furthermore, we will be able to determine the relationship between the variance of the growth rate in total factor productivity and the variance of the growth rate in output. Instead of using data across countries, we will use this data across states, thus reducing the variation in institutions, legal system, and tax rates.

We envision estimating the value of educational quality. By examining the effects of class size, teacher salary and student behavior on the return to schooling. Lazear (2001) points out that class size is determined by the public good aspects of classroom education. If one student misbehaves or asks a question other students do not share, the teacher must devote time to discipline or educate the lone interrupting student. Class size is a choice variable inversely related to the amount of disruption, and teacher salary. As students become better behaved, the efficient class size increases. Applying these data we can estimate the probability a student will behave in the class room by analyzing the joint teacher-salary/class size equilibrium. In so doing we believe the data will provide information on the returns to school quality, and add to the work of Card and Krueger (1992), Tamura (2001), Welch (1966) and others.

## REFERENCES

- Angrist, Joshua D. and Krueger, Alan B. "Does Compulsory School Attendance Affect Schooling and Earnings?", *Quarterly Journal of Economics*, 106, 1991: 979-1014.
- Arellano, M., and Bond, S. "Some Tests of Specification for Panel Data: Monte Carlo Evidence and an Application to Employment Equations," *Review of Economic Studies* 58, 1991: 277-297.
- Baier, Scott, Dwyer, Gerald, and Tamura, Robert. "How Important Are Capital and Total Factor Productivity for Economic Growth?" Clemson University working paper, 2004.
- Barro, Robert and Lee, Jong-Wha. "International Comparisons of Educational Attainment," *Journal of Monetary Economics* 32, 1993: 363-394.
- Berry, William D., Fording, Richard C., and Hanson, Russell L. "An Annual Cost of Living Index for the American States, 1960-95," *Journal of Politics*, 60 (May), 2000: 550-67.
- Bills, Mark and Klenow, Peter J. "Does Schooling Cause Growth?" *American Economic Review*. 90, December, 2000: 1160-1183.
- Card, David. "Using Geographic Variation in College Proximity to Estimate the Return to Schooling", in Louis N. Christofides, E. Kenneth Grant and Robert Swidinsky, eds. *Aspects of Labour Market Behaviour: Essays in Honour of John Vanderkamp*. University of Toronto Press, Toronto, Canada. 1995. 201-222.
- Card, David, and Krueger, Alan. "Does School Quality Matter? Returns to Education and the Characteristics of Public Schools in the United States," *Journal of Political Economy* 100, 1992: 1-40.
- Denison, Edward. *The Sources of Economic Growth in the United States and the Alternatives Before Us*, Committee for Economic Development: New York, 1962.
- Easterlin, Richard. "Regional Growth of Income: Long Term Tendencies, 1880-1950." In *Population Redistribution and Economic Growth, United States, 1870-1950*, vol. 1, *Analyses of Economic Change*, edited by Simon Kuznets, Ann Ratner Miller, and Richard Easterlin. Philadelphia: American Philosophical Society, 1960a.

- Easterlin, Richard. "Interregional Differences in Per Capita Income, Population, and Total Income, 1840-1950," in *Trends in the American Economy in the Nineteenth Century*, ed. William N. Parker, Princeton University Press: Princeton, 1960b.
- Freeman, Richard. *The Overeducated American*, Academic Press: New York, 1976.
- Goldin, Claudia. "America's Graduation from High School: The Evolution and Spread of Secondary Schooling in the Twentieth Century," *Journal of Economic History* 58, 1999: 345-374.
- Goldin, Claudia, and Margo, Robert A. "Wages, Prices, and Labor Markets before the Civil War," in *Strategic Factors in Nineteenth Century American Economic History: A Volume to Honor Robert W. Fogel* (eds.) Claudia Goldin and Hugh Rockoff, NBER, University of Chicago Press: Chicago, 1992a.
- Goldin, Claudia, and Margo, Robert A. "The Great Compression: The Wage Structure in the United States at Mid-Century," *Quarterly Journal of Economics* 107, 1992b: 1-34.
- Goldin, Claudia, and Katz, Lawrence. "Education and Income in the Early 20<sup>th</sup> Century: Evidence from the Prairies," *Journal of Economic History* 60, 2000: 782-818.
- Gordon, Robert. *Macroeconomics*, Addison-Wesley: New York, 1999.
- Kuznets, Simon. *National Income: A Summary of Findings*, National Bureau of Economic Research: New York, 1946.
- Lazear, Edward. "Educational Production" *Quarterly Journal of Economics*. 116, 2001:777-803
- Mitchener, Kris J. and McLean, Ian W. "U.S. Regional Growth and Convergence 1880 -1980." *Journal of Economic History*, 59, 1997: 1016-1042.
- Mulligan, Casey, and Sala-i-Martin, Xavier. "A Labor-Income Based Measure of the Aggregate Value of Human Capital," *Journal of Japan and the World Economy* 9, 1997: 159-191.
- Mulligan, Casey, and Sala-i-Martin, Xavier. "Measuring Aggregate Human Capital," *Journal of Economic Growth* 5, 2000: 215-252.

- Murphy, Kevin M., and Welch, Finis. "The Structure of Wages," *Quarterly Journal of Economics* 107, 1992: 285-326.
- National Catholic Education Association. *United States Catholic elementary and secondary schools*, National Catholic Education Association: Washington, D.C., various years.
- Snyder, Thomas, Hoffman, Leff, and Geddes, Claire. *State Comparisons of Education Statistics: 1969-70 to 1996-97*, National Center for Education Statistics: U.S. Department of Education: Washington, D.C., 1998.
- Staiger, Douglas and Stock, James H. "Instrumental Variables Regression with Weak Instruments", *Econometrica* 65, 1997: 557-586.
- Tamura, Robert. "Teachers, Growth and Convergence," *Journal of Political Economy* 109, 2001: 1021-1059.
- U.S. Census Bureau. *Statistical Abstracts of the United States*, U.S. Government Printing Office: Washington, D.C., various years.
- U.S. Department of Commerce. *Historical Statistics of the United States: Colonial Times to 1970*, U.S. Government Printing Office: Washington, D.C., 1975.
- U.S. Department of Education. *Digest of Education Statistics*, U.S. Government Printing Office: Washington, D.C., various years.
- U.S. Department of Health, Education and Welfare. *Projections of Educational Statistics to ...*, U.S. Government Printing Office: Washington, D.C., various years.
- Welch, Finis. "Measurement of the Quality of Schooling," *American Economic Review* 56, 1966: 379-392.
- Williamson, Jeffrey G. and Linder, Peter H. *American Inequality*, Academic Press 1980: 97-132.
- Wooldridge, Jeffrey. *Econometric Analysis of Cross Section and Panel Data*, M.I.T. Press: Cambridge, MA 2002.

## APPENDIX A

There are nine census regions in the US. The following Table provides the regional groups.

<i>New England</i>	<i>Middle Atlantic</i>	<i>South Atlantic</i>	<i>E. South Central</i>	<i>W. South Central</i>
Connecticut	New Jersey	Delaware	Alabama	Arkansas
Maine	New York	D.C.	Kentucky	Louisiana
Massachusetts	Pennsylvania	Florida	Mississippi	Oklahoma
New Hampshire		Georgia	Tennessee	Texas
Rhode Island		Maryland		
Vermont		North Carolina		
		South Carolina		
		Virginia		
		West Virginia		
<i>Mountain</i>	<i>Pacific</i>	<i>W. North Central</i>	<i>E. North Central</i>	
Arizona	Alaska	Iowa	Illinois	
Colorado	California	Kansas	Indiana	
Idaho	Hawaii	Minnesota	Michigan	
Montana	Oregon	Missouri	Ohio	
Nevada	Washington	Nebraska	Wisconsin	
New Mexico		North Dakota		
Utah		South Dakota		
Wyoming				

## APPENDIX B

In this Appendix we provide more details on the calculations of years of schooling.

- I. Describe collection/data
  - A. Public Elementary / Secondary Enrollment
  - B. Private Elementary / Secondary Enrollment
  - C. Higher Educational Enrollment
  - D. Population – 5-13, 14-17, 18-24, 65+. . . .
  - E. Labor Force
  - F. Price levels

- G. Expected years
- II. Describe calculation of
  - A. Enrollment rates
  - B. Shares (primary, secondary, college)
    - 1. Initial conditions
    - 2. Higher education inflow constant
    - 3. Secondary departure rates
  - C. Shares for foreign born
- III. Idiosyncrasies
  - A. DC / MD / VA
  - B. AK / HA
  - C. ND / SD / Dakota
  - D. OK / Indian Territory
- IV. Table listing first year of data availability
- V. References

## Data Description

### Public Enrollment Data.—

**Public Enrollment, 1840-1916** Data for total (elementary and secondary) public enrollment are available from decennial census data, by state, in 1840, 1850, 1860, 1870. Total public enrollment data are available in *Statistical Abstracts of the United States* for the years 1872, 1877, 1879-1887, 1889-1891, 1893-1916.

Data for total public enrollment for non-decennial years between 1840 and 1870 was geometrically interpolated. Data for the years 1871, 1873-1876, 1878, 1888, and 1892 was also geometrically interpolated.

We do not observe the fraction of total public enrollment that is elementary versus secondary until the year 1899. However, we do have national aggregates that make this breakdown in 1870, 1880, and 1890-1898.

Letting  $pub.enroll_{it}^{primary}$  designate the public primary enrollment level in state  $i$  for time period  $t$ ,



and  $pub.enroll_{it}^{total}$  refer to the total (primary and secondary) enrollment level, we assign:

$$pub.enroll_{it}^{primary} = pub.enroll_{it}^{total} \frac{\sum_j pub.enroll_{j,1870}^{primary}}{\sum_j pub.enroll_{j,1870}^{total}}, t \leq 1870 \quad (27)$$

$$pub.enroll_{it}^{primary} = pub.enroll_{it}^{total} \frac{\sum_j pub.enroll_{j,1880}^{primary}}{\sum_j pub.enroll_{j,1880}^{total}}, 1871 \leq t \leq 1880 \quad (28)$$

$$pub.enroll_{it}^{primary} = pub.enroll_{it}^{total} \frac{\sum_j pub.enroll_{j,1890}^{primary}}{\sum_j pub.enroll_{j,1890}^{total}}, 1881 \leq t \leq 1890 \quad (29)$$

$$pub.enroll_{it}^{primary} = pub.enroll_{it}^{total} \frac{\sum_j pub.enroll_{jt}^{primary}}{\sum_j pub.enroll_{jt}^{total}}, 1891 \leq t \leq 1898 \quad (30)$$

$$pub.enroll_{it}^{secondary} = pub.enroll_{it}^{total} - pub.enroll_{it}^{primary} \quad (31)$$

Beginning in 1899, we observe both  $pub.enroll_{it}^{total}$  and  $pub.enroll_{it}^{secondary}$  so we can simply calculate  $pub.enroll_{it}^{primary}$ .

**Public Enrollment, 1918 - 1968** Data for public secondary enrollment and for total public enrollment are available biennially in the *Statistical Abstract of the United States* (even numbered years) from 1918–1968. In addition, data is also available in 1925, 1945, 1947, and 1949, 1955, and 1959. We geometrically interpolate any missing values from 1918–1968.

**Public Enrollment, 1969 - 2000** Data from 1969 to 2000 are annual, and come from NCES, *State Comparisons of Education Statistics: 1969-70 to 1996-97*, as well as updates available from the NCES website.

#### Private Enrollment Data.—

**Private Enrollment, 1840 - 1916** Data for total private enrollments are available from various censuses, by state in 1840, 1850, 1860, 1870, 1890, 1910, and 1920. We geometrically interpolate between the decennial values listed above for any non-decennial years.

Data for total private secondary enrollments are available on an annual basis from 1899 to 1916 from the *Statistical Abstracts of the United States*. For these years, we are able to take the measure

of total private enrollment above and subtract secondary enrollment to arrive at private elementary enrollment.

Prior to 1899, we observe total private enrollment, but do not observe the breakdown into elementary and secondary. However, we do observe national aggregates in 1890. Proceeding as we did above in the public case, we calculate:

$$pri.enroll_{it}^{primary} = pri.enroll_{it}^{total} \frac{\sum_j pri.enroll_{j,1890}^{primary}}{\sum_j pri.enroll_{j,1890}^{total}}, t \leq 1890 \quad (32)$$

$$pri.enroll_{it}^{secondary} = pri.enroll_{it}^{total} - pri.enroll_{it}^{primary} \quad (33)$$

We also geometrically interpolate the secondary enrollment figures for 1891-1898 using the 1890 value (calculated directly above), and the 1899 figures.

**Private Enrollment, 1918 - 1968** Data for private secondary enrollment and total private enrollment are available biennially in *Statistical Abstracts of the United States* (even numbered years) from 1918–1940 and 1948–1968. Data is also available in 1925, 1947, and 1949, 1955, and 1959. We geometrically interpolate any missing values from 1918 – 1968.

**Private Enrollment, 1969 - 2000** For the years 1968 – 1980, 1991, 1993, 1995, 1997, and 1999, we observe private elementary and secondary enrollment figures from the *Digest of Education Statistics*. We geometrically interpolate the 1992, 1994, 1996, and 1998 values.

For the years between 1980 through 1991, we are unable to obtain private elementary and private secondary enrollment figures by state directly. However we are able to obtain annual estimates of the national private elementary and private secondary totals from Projections of Education Statistics, various issues, as well as state level data on Catholic elementary and Catholic secondary enrollment figures in 1985, 1988, and 1990 – 1999 from the *National Catholic Education Association*, various issues. We assume that the distribution of total private elementary and total private secondary enrollment figures across states is identical to the distribution of Catholic elementary and Catholic secondary enrollment figures across states. We inflate the Catholic state level data enrollment data to correspond to the national totals for 1985, 1988, and 1990. We geometrically interpolate values for years 1981-1984, 1986-1987, and 1988.

**Higher Education Enrollment.—**

#### *1840 – 1899*

Data for states are available from decennial census data in 1840, 1850, 1860, and 1870. In 1886, 1890, and 1891 data are available, typically subdivided into Medical, Theological, Law, and Liberal Arts enrollments. Data for non-census years between 1840 and 1870, as well as 1871-1885, 1887-1889, and 1892-1898 are geometrically interpolated.

#### *1899 – 1920*

Data are reported annually in *Statistical Abstracts* under a variety of titles and formats. Total higher education enrollment is the sum of sources below, except where enrollment figures are included in more than one source.

1. Schools of Technology and Institutions conferring only the B.S. degree (1899-1905)
2. Colleges and Seminaries for Women which confer degrees (1899-1910)
3. Coeducational Colleges and Universities and Colleges for men only (1899-1916, 1918)
4. Undergraduate Students in Univ., Colleges, and Schools of Tech. (1911 – 1916, 1918, 1920)
5. Professional Schools (1899-1916)
6. Public and Private Normal Schools (1899-1916, 1918, 1920)
7. Training Schools for Nurses, Comm. Schools, Manual and Industrial Training Schools (1910-1916, 1918, 1920)

#### *1922 – 1946*

Data is reported biennially in the *Statistical Abstracts* from 1922-1940, various issues, as Enrollment in Universities, Colleges, and Preparatory Schools. Similar data is also reported as Higher Education Enrollment in 1942, 1944, and 1946. Non-biennial years are geometrically interpolated.

#### *1947 – 1968*

Data is reported annually in *Statistical Abstracts*, various issues, as Institutions of Higher Educational, Fall Enrollment.

#### *1969 – 2000*

Data is reported in *State Comparisons of Education Statistics*. Higher educational enrollment is the sum of 2-year private, 2-year public, 4-year private, and 4-year public higher educational enrollment.

### **Population.—**

We generally observe the age distribution of population in decennial years, beginning in 1840. In most cases, we are given data with 5-year population distributions. The usual structure is

<5, 5-9, 10-14, 15-19, 20-24... 55-59, 60-64, 65-69, 70-74. . .

With the exception of calculating the average age of the population in a state, we are ultimately interested in the age groups: 5-13, 14-17, 18-24, 16-65. In order to calculate the number of persons in each group, we assume a uniform distribution of population across the age groups.

In 1840, the white age distribution is reported, but only broad categories of the black age distribution are available. In order to allocate the total black distribution amongst the various age groups, we assume the fraction of total black population in each age group is identical to the fraction in the 1850 black distribution.

#### **Labor Force.**—

All labor force data prior to 1970 is decennial data. For non-decennial years prior to 1970, data is geometrically interpolated. Labor force data for 1840 – 1860 is decennial census data. Data for 1870 – 1940 is gainful workers, 10 years old and over, and is taken from *Historical Statistics of the United States: Colonial Times to 1970*, pp. 129–131. Data for 1950 and 1960 is decennial Census of Population data, and includes persons aged 14 and over. Data from 1970 – 2000 is Civilian Labor Force, 16 years and older, and is taken from the Bureau of Labor Statistics website.

#### **Price Levels.**—

National price level data from 1875-1999 is the GDP deflator, as reported in Gordon, *Macroeconomics*, 7th edition, pp. A1–A3. National price level data prior to 1875 is the wholesale price index (all commodities) from Warren and Pearson, printed in *Historical Statistics of the United States: Colonial Times to 1970*, pp. 201-202. Data from 1840-1875 are normalized to correspond to the price level given by Gordon in 1875.

In addition, we use three sources of information on relative price levels across regions. Mitchener and McLean (1999) and Williamson and Linder (1980) provide regional price levels for census regions which we use from 1840-1960. Data from the two sources is primarily non-overlapping. Where we have data from both sources, we take the arithmetic average of the relative price level in each region. Prior to 1880 these sources does not include relative price levels for the Pacific and Mountain region. For data prior to 1880 in each of these two regions, we *extrapolate* the relative regional price level using the trend observed from 1880 to 1920. Berry, Fording and Hanson (2000) display price levels for each state on an annual basis from 1960-2000. To maintain consistency, we aggregate these state level estimates into census regions. In non-decennial years, we interpolate relative price levels. We

normalize regional price levels in all years to the national price level figures given in Gordon (and Warren and Pearson). All income measures are reported in 2000 dollars.

**Expected Years.—**

The portion of the population, 25 years old and over that has completed various levels of school is given in the Census of the Population in 1940 – 2000. From this information, we calculate the expected number of years of school completed, conditional on being in either the primary, secondary, or higher educational group. The values for  $yr_s_t^{\text{college}}$ ,  $yr_s_t^{\text{secondary}}$ , and  $yr_s_t^{\text{primary}}$  were obtained from decennial census data. Let  $N(i - j)$  be the number of people who have completed between  $i$

and  $j$  years of schooling, inclusive.<sup>24</sup>

$$yrs_{1940}^{\text{primary}} = \frac{2.5N(1-4) + 5.5N(5-6) + 7.5N(7-8)}{N(1-4) + N(5-6) + N(7-8)} \quad (34)$$

$$yrs_{1950,1960,1970,1980}^{\text{primary}} = \frac{2.5N(1-4) + 5.5N(5-6) + 7N(7) + 8N(8)}{N(1-4) + N(5-6) + N(7) + N(8)} \quad (35)$$

$$yrs_{1990}^{\text{primary}} = \frac{2.5N(1-4) + 7.23N(5-8)}{N(1-4) + N(5-8)} \quad (36)$$

$$yrs_{2000}^{\text{primary}} = \frac{6.46N(0-8)}{N(0-8)} \quad (37)$$

$$yrs_{1940}^{\text{secondary}} = \frac{10N(9-11) + 12N(12)}{N(9-11) + N(12)} \quad (38)$$

$$yrs_{1950,1960,1970}^{\text{secondary}} = \frac{10N(9-11) + 12N(12)}{N(9-11) + N(12)} \quad (39)$$

$$yrs_{1980}^{\text{secondary}} = \frac{9N(9) + 10N(10) + 11N(11) + 12N(12)}{N(9) + N(10) + N(11) + N(12)} \quad (40)$$

$$yrs_{1990,2000}^{\text{secondary}} = \frac{10.5N(9-12) + 12N(12)}{N(9-12) + N(12)} \quad (41)$$

$$yrs_{1940,1950,1960}^{\text{college}} = \frac{14N(13-15) + 17N(16^+)}{N(13-15) + N(16^+)} \quad (42)$$

$$yrs_{1970}^{\text{college}} = \frac{14N(13-15) + 16N(16) + 18N(17^+)}{N(13-15) + N(16) + N(17^+)} \quad (43)$$

$$yrs_{1980}^{\text{college}} = \frac{13N(13) + 14N(14) + 15N(15) + 16N(16) + 17.5N(17-18) + 20N(19^+)}{N(13) + N(14) + N(15) + N(16) + N(17-18) + N(19^+)} \quad (44)$$

$$yrs_{1980}^{\text{college}} = \frac{13N(13) + 14N(14) + 15N(15) + 16N(16) + 17.5N(17-18) + 20N(19^+)}{N(13) + N(14) + N(15) + N(16) + N(17-18) + N(19^+)} \quad (45)$$

---

<sup>24</sup>This is the population 25 and older, not labor force, because we are looking at only those who have presumably completed schooling.

$$yrs_{1990}^{\text{college}} = \frac{14N(\text{scn} + a) + 16N(b) + 18N(\text{ma}) + 19.75N(\text{pr}) + 20N(d)}{N(\text{scn}) + N(a) + N(b) + N(\text{ma}) + N(\text{pr}) + N(d)} \quad (52)$$

(53)

$$yrs_{2000}^{\text{college}} = \frac{14N(\text{sc}) + 14N(a) + 16N(b) + 18N(\text{ma}) + 19.75N(\text{prg})}{N(\text{sc}) + N(a) + N(b) + N(\text{ma}) + N(\text{prg})} \quad (54)$$

9 – 12 = 9th to 12th grade, no diploma

*sc* = some college

*scn* = some college no degree

*a* = Associate degree

*b* = Bachelor's degree

*ma* = Master's degree

*prg* = Professional. or Graduate degree

*pr* = Professional school degree

*d* = Doctorate degree

In 1990, data are not reported as finely for those who have completed between 5 and 8 years of schooling. We need to assign a number of years of schooling to give to the group  $N(5 - 8)$ , but this distribution is highly skewed. We calculate the conditional distribution in the years 1960, 1970, and 1980. We assign 7.23 years in 1990.

$$yrs_{1960}^{5-8} = 5.5N(5 - 6)_{1960} + 7N(7)_{1970} + 8N(8)_{1960} = 7.22 \quad (55)$$

$$yrs_{1970}^{5-8} = 5.5N(5 - 6)_{1970} + 7N(7)_{1970} + 8N(8)_{1970} = 7.23 \quad (56)$$

$$yrs_{1980}^{5-8} = 5.5N(5 - 6)_{1980} + 7N(7)_{1980} + 8N(8)_{1980} = 7.24 \quad (57)$$

$$yrs_{1990}^{5-8} = 7.23 \quad (58)$$

In 2000, we need to assign a number of years of schooling to give to the group  $N(0 - 8)$ , whose distribution is highly skewed. We use March 2000 CPS data for the population of people age 15 or over, which gives us data that is less aggregated than the census data. We assign 7.74 years to  $N(7 - 8)$ , which is the average value from the 1960 (7.73), 1970 (7.75), and 1980 (7.75)  $yrs^{5-8}$ . Thus the calculated value for  $yrs_{2000}^{0-8}$  is 6.42:

$$yrs_{2000}^{0-8} = 2.5N(1 - 4)_{2000} + 5.5N(5 - 6)_{2000} + 7.74N(7 - 8)_{2000} = 6.42 \quad (59)$$

Values for  $yr_s^i$  for periods prior to 1940 were calculated by geometrically interpolating from an initial value for the year in which the state first has adequate data available (see Table A1) to the 1940 value. Initial values are 4, 10, and 14 for primary, secondary, and higher education, respectively.

All values for non-census years between 1940 and 2000 were geometrically interpolated. We do not include those persons for whom the educational attainment level is not reported.

## Description of Calculations

### General Enrollment Rates.—

Enrollment figures for public and private school are summed to obtain a total primary enrollment rate, total secondary enrollment rate, and total higher educational enrollment rate. From enrollment data, enrollment rates are calculated as below:

$$tot.enroll_t^{\text{primary}} = pub.enroll_t^{\text{primary}} + pri.enroll_t^{\text{primary}} \quad (60)$$

$$tot.enroll_t^{\text{secondary}} = pub.enroll_t^{\text{secondary}} + pri.enroll_t^{\text{secondary}} \quad (61)$$

$$tot.enroll_t^{\text{college}} = pub.enroll_t^{\text{college}} + pri.enroll_t^{\text{college}} \quad (62)$$

$$r_t^{\text{primary}} = \frac{tot.enroll_t^{\text{primary}}}{\ell[5 - 13]_t} \quad (63)$$

$$r_t^{\text{secondary}} = \frac{tot.enroll_t^{\text{secondary}}}{\ell[14 - 17]_t} \quad (64)$$

$$r_t^{\text{college}} = \frac{tot.enroll_t^{\text{college}}}{\ell[18 - 24]_t} \quad (65)$$

### General educational exposure shares.—

To calculate the stock of human capital of each type, primary school stock, secondary school stock and higher education stock, we used a perpetual inventory method. The following will illustrate the nature of our calculations. We ignore state subscripts without loss of information. In period  $t+1$ , the stock of adults, with exposure to education level  $i$ ,  $i$ =primary, secondary, and higher, but no more is given by:

$$H_{t+1}^i = H_t^i(1 - \delta_t^i) + I_t^i \quad (66)$$

where  $\delta_t^i$  is the death rate and  $I_t^i$  is the flow of new adults with exposure to education level  $i$  and no more. Initially, we assume that  $\delta_t$  does not vary by education class, while later we estimate the death rate separately for the secondary and higher educational classes.

It is useful to put the human capital measure as a fraction of the labor force. Thus, we normalize



and produce

$$\frac{H_{t+1}^i}{L_{t+1}} = \frac{H_t^i}{L_t} \frac{L_t}{L_{t+1}} (1 - \delta_t) + \frac{I_t^i}{L_{t+1}} \quad (67)$$

$$h_{t+1}^i = h_t^i \frac{L_t}{L_{t+1}} (1 - \delta_t) + \frac{I_t^i}{L_{t+1}} \quad (68)$$

where  $h_t^i$  measures the share of the labor force exposed to education level  $i$ , and no more in year  $t$ .

The flows into education categories are given by:

$$I_t^{\text{college}} = \frac{r_t^{\text{college}} l f p r_t^{\text{college}} \ell[18-24]_t}{7} \Theta \quad (69)$$

$$I_t^{\text{secondary}} = \frac{(r_t^{\text{secondary}} - r_t^{\text{college}} \Theta) l f p r_t^{\text{secondary}} \ell[14-17]_t}{4} \quad (70)$$

$$I_t^{\text{primary}} = \frac{(r_t^{\text{primary}} - r_t^{\text{secondary}}) l f p r_t^{\text{primary}} \ell[5-13]_t}{9} \quad (71)$$

$$I_t^{\text{none}} = \frac{(1 - r_t^{\text{primary}}) l f p r_t^{\text{none}} \ell[5-13]_t}{9} \quad (72)$$

where  $r_t^i$   $i$ =college, secondary and primary are the respective enrollment rates,  $l f p r_t^i$  are the labor force participation rates for education category  $i$  including those without schooling,  $\ell[i-j]$  is the number of people between the ages of  $i$  and  $j$ , inclusive, and  $\Theta$  is the constant to adjust the inflow into the higher educational category.

In order to proceed we need a measure of  $\delta_t^i$ , the death rate of adults.

As  $L_{t+1} = L_t (1 - \delta_t) + I_t^{\text{college}} + I_t^{\text{secondary}} + I_t^{\text{primary}} + I_t^{\text{none}}$

Using the above definitions, notice that this allows for the calculation of  $\frac{L_t}{L_{t+1}} (1 - \delta_t)$ :

$$\frac{L_t}{L_{t+1}} (1 - \delta_t) = 1 - \frac{\frac{r_t^{\text{college}} l f p r_t^{\text{college}} \ell[18-24]_t}{7} \Theta + \frac{(r_t^{\text{secondary}} - r_t^{\text{college}} \Theta) l f p r_t^{\text{secondary}} \ell[14-17]_t}{4} + \frac{(r_t^{\text{primary}} - r_t^{\text{secondary}}) l f p r_t^{\text{primary}} \ell[5-13]_t}{9} + \frac{(1 - r_t^{\text{primary}}) l f p r_t^{\text{none}} \ell[5-13]_t}{9}}{L_{t+1}} \quad (73)$$

With this information, we can calculate each of the shares of the labor force with each schooling category.

Using this method produced a much smaller share of the labor force exposed to higher education than the census figures. Thus we estimate the death rate of those exposed to higher education independently. We assumed that there was no death, just retirement from the labor force after 45 years of work. We have the stock of adults exposed to higher education, given as follows:

$$H_{t+1}^{\text{college}} = H_t^{\text{college}} - I_{t-45}^{\text{college}} + I_t^{\text{college}} \quad (74)$$

$$\frac{H_{t+1}^{\text{college}}}{L_{t+1}} = \frac{H_t^{\text{college}}}{L_t} \frac{L_t}{L_{t+1}} - \frac{I_{t-45}^{\text{college}}}{L_{t-45}} \frac{L_{t-45}}{L_{t+1}} + \frac{I_t^{\text{college}}}{L_{t+1}} \quad (75)$$

Thus, to calculate the higher education share in period  $t$ , we must measure  $\frac{I_{t-45}^{\text{college}}}{L_{t-45}}$ , which requires higher education enrollment data in period  $t-45$ . For the earlier portion of our sample, we do not observe enrollment rates early enough to make this calculation. Where necessary, we linearly interpolate between the 0 and the value of the higher education enrollment rate the first time it is observed. See Table B.2 for the years in which each state is first calculated, and for the first time we observe higher educational enrollment figures. Unfortunately we do not observe  $L_{t-45}$  until we have 45 years of state data. We use the labor force participation rate closest to that year and then use population data to calculate a measure of the labor force  $L_{t-45}$ .

Initially, we had included the secondary education exposed workers in a category along with those workers exposed to elementary education and no education. However, we find that this resulted in calculated shares exposed to elementary education that were less than zero. As a result, we choose  $\delta^{\text{secondary}}$  for each state by matching the calculated shares of workers exposed to secondary education to those observed in the census years from 1940-2000. This allows us to calculate the share of workers exposed to secondary education:

$$h_{t+1}^{\text{secondary}} = h_t^{\text{secondary}} \frac{L_t}{L_{t+1}} (1 - \delta^{\text{secondary}}) + \frac{I_t^{\text{secondary}}}{L_{t+1}} \quad (76)$$

Given that we have calculated for  $h_t^{\text{college}}$  and  $h_t^{\text{secondary}}$  in all periods, we can proceed to calculate the shares for primary and no schooling. The next set of equations shows how we can identify the term  $\frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}})$ .

$$L_{t+1} = H_{t+1}^{\text{college}} + H_{t+1}^{\text{secondary}} + H_{t+1}^{\text{primary}} + H_{t+1}^{\text{none}} \quad (77)$$

$$L_{t+1} = H_t^{\text{college}} - I_{t-45}^{\text{college}} + I_t^{\text{college}} + H_t^{\text{secondary}} (1 - \delta^{\text{secondary}}) + (H_t^{\text{primary}} + H_t^{\text{none}}) (1 - \delta_t^{\text{primary}}) + (I_t^{\text{secondary}} + I_t^{\text{primary}} + I_t^{\text{none}}) \quad (78)$$

$$1 - h_{t+1}^{\text{college}} = h_t^{\text{secondary}} \frac{L_t}{L_{t+1}} (1 - \delta^{\text{secondary}}) + (h_t^{\text{primary}} + h_t^{\text{none}}) \frac{L_t}{L_{t+1}} (1 - \delta_t^{\text{primary}}) + \frac{I_t^{\text{secondary}} + I_t^{\text{primary}} + I_t^{\text{none}}}{L_{t+1}} \quad (79)$$

$$\frac{L_t}{L_{t+1}} \left(1 - \delta_t^{\text{primary}}\right) = \frac{1 - h_{t+1}^{\text{college}} - h_{t+1}^{\text{secondary}} \frac{L_t}{L_{t+1}} \left(1 - \delta^{\text{secondary}}\right) - \left(\frac{I_t^{\text{secondary}}}{L_{t+1}} + \frac{I_t^{\text{primary}}}{L_{t+1}} + \frac{I_t^{\text{none}}}{L_{t+1}}\right)}{\left(h_t^{\text{primary}} + h_t^{\text{none}}\right)} \quad (80)$$

$$\frac{L_t}{L_{t+1}} \left(1 - \delta_t^{\text{primary}}\right) = \frac{1 - h_{t+1}^{\text{college}} - h_{t+1}^{\text{secondary}} \frac{L_t}{L_{t+1}} \left(1 - \delta^{\text{secondary}}\right)}{\left(h_t^{\text{primary}} + h_t^{\text{none}}\right)} - \frac{\left(\frac{\left(r_t^{\text{secondary}} - r_t^{\text{college}} \Theta\right) l f p r_t^{\text{secondary}} \ell[14-17]_t}{4} + \frac{\left(r_t^{\text{primary}} - r_t^{\text{secondary}}\right) l f p r_t^{\text{primary}} \ell[5-13]_t}{9} + \frac{\left(1 - r_t^{\text{primary}}\right) l f p r_t^{\text{none}} \ell[5-13]_t}{9}\right)}{L_{t+1} \left(h_t^{\text{primary}} + h_t^{\text{none}}\right)} \quad (81)$$

We occasionally measure primary and secondary enrollment rates that are larger than unity. There are a couple of reasons why this occurs. The data contains individuals that were held back in school, and also there are people that receive education for the first time starting at an unusual age. Since we have very limited information on repeaters as well as unusual starters, we treat all cases as the latter.

**Initial Conditions** The initial condition for  $h_t^i$ ,  $i = \text{college, secondary and primary}$  were the respective enrollment rate of each class divided by two.

### Higher Ed Inflow Adjustment & Secondary Departure Rates

Table B1. Values of  $\Theta$  and  $\delta^{\text{secondary}}$

New England	$\Theta$	$\delta^{\text{secondary}}$	E. South Central	$\Theta$	$\delta^{\text{secondary}}$	W. North Central	$\Theta$	$\delta^{\text{secondary}}$
Connecticut	1.37	.990	Alabama	1.24	.982	Iowa	1.17	.976
Maine	1.27	.977	Kentucky	1.17	.966	Kansas	1.23	.972
Massachusetts	1.18	.979	Mississippi	1.21	.971	Minnesota	1.30	.978
New Hampshire	1.27	.986	Tennessee	1.18	.977	Missouri	1.15	.975
Rhode Island	1.15	.978				Nebraska	1.15	.968
Vermont	1.24	.978				North Dakota	1.19	.943
						South Dakota	1.19	.959

Middle Atlantic	W. South Central				E. North Central			
New Jersey	1.31	.989	Arkansas	1.21	.970	Illinois	1.27	.988
New York	1.27	.983	Louisiana	1.21	.970	Indiana	1.18	.981
Pennsylvania	1.00	.973	Oklahoma	1.23	.968	Michigan	1.26	.988
			Texas	1.37	.990	Ohio	1.00	.981
						Wisconsin	1.27	.976

South Atlantic	Mountain				Pacific			
Delaware	1.36	.998	Arizona	1.45	.999	Alaska	1.65	.999
D.C.	1.25	.973	Colorado	1.60	.998	California	1.27	.9995
Florida	1.27	.999	Idaho	1.51	.990	Hawaii	1.52	.984
Georgia	1.48	.992	Montana	1.48	.973	Oregon	1.35	.999
Maryland	1.27	.9995	Nevada	1.45	.9999	Washington	1.36	.998
North Carolina	1.27	.980	New Mexico	1.35	.991			
South Carolina	1.27	.981	Utah	1.36	.983			
Virginia	1.30	.990	Wyoming	1.48	.989			
West Virginia	1.00	.968						

**Foreign Shares.**—

In the calculation of our measure of years of schooling in state  $i$ , recall that we multiply the fraction of state  $i$ 's residents that were born in state  $j$  by the years of schooling in state  $j$  (assuming

no mobility):

$$E_{it} = \sum_{j \neq for} S_{ijt} \widehat{E}_{jt} \quad (82)$$

We derived our measure of  $\widehat{E}_{jt}$  from observing the enrollment rates in state  $j$  and using the perpetual inventory methodology described above. Because a fraction of the residents of state  $i$ 's residents are foreign born, we require a measure of  $\widehat{E}_{for,t}$ , the average years of schooling for the foreign born. If we could observe the share of the foreign born in each education category, we would simply calculate:

$$\widehat{E}_{for,t} = h_{for,t}^{\text{primary}} yrs_{for,t}^{\text{primary}} + h_{for,t}^{\text{secondary}} yrs_{for,t}^{\text{secondary}} + h_{for,t}^{\text{college}} yrs_{for,t}^{\text{college}} \quad (83)$$

However, this data is not available, and thus we cannot calculate the corresponding measures of  $h_{for,t}^{\text{primary}}$ ,  $h_{for,t}^{\text{secondary}}$  and  $h_{for,t}^{\text{college}}$ .

We use two different adjustment algorithms. We initially calculate the average years of schooling excluding the contributions made by the foreign born, which we denote  $\widetilde{E}_{it}$ :

$$\widetilde{E}_{it} = \sum_{j \neq for} S_{ijt} \widehat{E}_{jt} \quad (84)$$

We then assign the number of years of schooling to the foreign born  $\widehat{E}_{for,t}$  so that our overall years of schooling measure,  $E_{it}$  equals the years of schooling reported by the census,  $yrs_{it}$ :

$$\widehat{E}_{for,t} = \frac{\left( yrs_{it} - \widetilde{E}_{it} \right)}{S_{i,for,t}} \quad (85)$$

We then place a lower and upper bound on average years of schooling assigned to foreigners by:

$$\widehat{E}_{for,t} \in \left[ 1, yrs_{it}^{\text{college}} \right] \quad (86)$$

We allocate the shares among the educational categories such that:

$$\widehat{E}_{for,t} = \widehat{h}_{for,t}^{\text{primary}} yrs_{it}^{\text{primary}} + \widehat{h}_{for,t}^{\text{secondary}} yrs_{it}^{\text{secondary}} + \widehat{h}_{for,t}^{\text{college}} yrs_{it}^{\text{college}} \quad (87)$$

Although there is no unique allocation, we assigned the shares using the following algorithm, in order to preserve the equality of (87):

If  $\widehat{E}_{for,t} < yrs_{it}^{\text{primary}}$ , we allocate between the none and primary categories, assigning zero for the secondary and college. In this case,  $\widehat{E}_{for,t} = \frac{yrs_{it}^{\text{primary}}}{S_{i,for,t}}$  and  $\widehat{h}_{for,t}^{\text{none}} = \left( 1 - \widehat{h}_{for,t}^{\text{primary}} \right)$ . If  $yrs_{it}^{\text{primary}} <$

$\widehat{E}_{for,t} < yrs_{it}^{secondary}$ , we assign zero for the none and college categories and allocate between the primary and secondary categories. If  $yrs_{it}^{secondary} < \widehat{E}_{for,t} < yrs_{it}^{college}$ , we assign zero for the none and primary categories and allocate between the secondary and college groups. If  $\widehat{E}_{for,t} > yrs_{it}^{college}$ , we allocate between the secondary and college categories, assigning zero for the none and primary.

**California Adjustment** The algorithm above assumes that the foreign born population is homogeneous in regards to educational attainment. If the number of years assigned to foreigners lies between  $yrs_{it}^{primary}$  and  $yrs_{it}^{secondary}$ , the algorithm would assign foreigners a zero share to the college and none categories. If the *actual* distribution of foreigners contains a substantial fraction of workers categorized as none and college, the algorithm would mistakenly assign these workers into the primary and secondary categories. While this is a possibility in all states, we feel this is particularly troublesome in California after 1970. In California it is quite plausible that the foreign born may be comprised of two distinct groups - a highly educated group, and a group of new migrants with low educational attainment levels. Using this algorithm for California, we would overestimate primary and secondary, but more importantly, underestimate college. This problem is further exacerbated by a growing share of the population that is foreign born. This would result in a substantial underestimation of the share exposed to college after 1970. To address this problem, we assign half of the foreign born to the college category after 1980.<sup>25</sup> We then allocate the remaining years to be assigned between secondary and primary. The remaining foreign born are assigned to the none category.

## Idiosyncrasies

### DC / MD / VA.—

We observe extremely high private enrollment rates for District of Columbia throughout the sample, presumably due to a large number of non-residents attending the District of Columbia schools. We surmise that these enrollment figures are overstated as many residents of Maryland and Virginia are attending District of Columbia schools.

From 1910 – 1999, we assign a private elementary enrollment rate equal to zero for DC. We apportion those private elementary students enrolled in DC into the private elementary enrollment

---

<sup>25</sup>We linearly interpolate the value of  $h_{it}^{college}$  between 1970 and 1980.

figures for Maryland and Virginia, using the population aged 5-13.

$$\begin{aligned}
pri.enroll_{Md,t}^{primary} &= pri.enroll_{Md,t}^{primary} \\
&+ \left( \frac{\ell[5-13]_{Md,t}}{\ell[5-13]_{Va,t} + \ell[5-13]_{Md,t}} \right) pri.enroll_{DC,t}^{primary} \quad (88)
\end{aligned}$$

$$\begin{aligned}
pri.enroll_{Va,t}^{primary} &= pri.enroll_{Va,t}^{primary} \\
&+ \left( \frac{\ell[5-13]_{Va,t}}{\ell[5-13]_{Va,t} + \ell[5-13]_{Md,t}} \right) pri.enroll_{DC,t}^{primary} \quad (89)
\end{aligned}$$

We allow the private secondary enrollment rate in DC to be no higher than the private secondary enrollment rate in the state of Massachusetts. We first calculate the enrollment rate in excess of the enrollment rate in DC, and then calculate the implied excess enrollment (students). We then apportion the excess enrollment into MD and VA, weighted by the population aged 14-17 in each state.

$$pri.enroll_{DC,t}^{secondary} = pri.r_{Ma,t}^{secondary} \ell[14-17]_{DC,t} \quad (90)$$

$$\begin{aligned}
pri.enroll_{Md,t}^{secondary} &= pri.enroll_{Md,t}^{secondary} \\
&+ \left( \frac{\ell[14-17]_{Md,t} \cdot (pri.r_{DC,t}^{secondary} - pri.r_{Ma,t}^{secondary})}{\ell[14-17]_{Va,t} + \ell[14-17]_{Md,t}} \right) \ell[14-17]_{DC,t} \quad (91)
\end{aligned}$$

$$\begin{aligned}
pri.enroll_{Va,t}^{secondary} &= pri.enroll_{Va,t}^{secondary} \\
&+ \left( \frac{\ell[14-17]_{Va,t} \cdot (pri.r_{DC,t}^{secondary} - pri.r_{Ma,t}^{secondary})}{\ell[14-17]_{Va,t} + \ell[14-17]_{Md,t}} \right) \ell[14-17]_{DC,t} \quad (92)
\end{aligned}$$

**AK / HA.**—

$yrst^{college}$ ,  $yrst^{secondary}$ , and  $yrst^{primary}$  for Alaska in 1939 and for Hawaii in 1940 were set as 14.5, 10.5, and 5.5 respectively.

**ND / SD/ Dakota.**—

From 1880 through 1890, population and enrollment figures are reported for Dakota, which is the aggregate of North Dakota and South Dakota. In 1890, we first observe separate figures for North Dakota and South Dakota. Where data is available, we allocate a constant fraction of Dakota population and enrollment figures to each of North and South Dakota, based on the population of each state in 1890.

**Indian Territory / Oklahoma.—**

We first include Oklahoma in our data set only after the *Statistical Abstract* reported data for Oklahoma, rather than Indian Territory.

Table B2: List of first year we observe enrollment data, and first year we observe higher education

enrollment data.					
State	1 <sup>st</sup> year of obs.	1 <sup>st</sup> year of higher ed.	State	1 <sup>st</sup> year of obs.	1 <sup>st</sup> year of higher ed.
Alabama	1840	1840	Montana	1870	1870
Alaska	1939	1924	Nebraska	1860	1870
Arizona	1872	1899	Nevada	1870	1886
Arkansas	1840	1850	New Hampshire	1840	1840
California	1850	1860	New Jersey	1840	1840
Colorado	1870	1870	New York	1840	1840
Delaware	1840	1840	North Carolina	1840	1840
D.C.	1850	1850	North Dakota	1890	1890
Florida	1840	1870	Ohio	1840	1840
Georgia	1840	1840	Oklahoma	1890	1899
Hawaii	1940	1922	Oregon	1850	1860
Idaho	1870	1899	Pennsylvania	1840	1840
Illinois	1840	1840	Rhode Island	1840	1840
Indiana	1840	1840	South Carolina	1840	1840
Iowa	1840	1850	South Dakota	1890	1890
Kansas	1860	1860	Tennessee	1840	1840
Kentucky	1840	1840	Texas	1850	1850
Louisiana	1840	1840	Utah	1860	1870
Maine	1840	1840	Vermont	1840	1840
Maryland	1840	1840	Virginia	1840	1840
Massachusetts	1840	1840	Washington	1860	1870
Michigan	1840	1840	West Virginia	1870	1870
Minnesota	1860	1860	Wisconsin	1850	1850
Mississippi	1840	1840	Wyoming	1870	1890
Missouri	1840	1840			



## APPENDIX C

To analyze the return to schooling, we need information on the income per worker. Since 1929, the Bureau of Economic Analysis has reported state level annual income data. Total and per capita state income for 1840, 1880, 1900 and 1919-1921 are documented by Richard Easterlin in his works, “Interregional Differences in Per Capita Income, Population, and Total Income 1840-1950” in *Trends in the American Economy in the Nineteenth Century* and *Analyses of Economic Change in Population Redistribution and Economic Growth, United States, 1870-1950*. These data exclude transfer payments, likely small during this time period, and the figures for 1840 do not include all components of personal income. For the Census years not reported by Easterlin, 1850, 1860, 1870, 1890, and 1910, we generate the missing state per capita income using data available from the Easterlin sources above, the 1850 through 1910 Censuses, and the *Historical Statistics of the United States: Colonial Times to 1970* (HSUS). In order to calculate state per worker income, we calculate value added by each industry at the state level. Although data is not available for every industry, production value is reported for agriculture in the Census from 1870 through 1910 and production value and materials are reported in the Census from 1850 through 1910 for manufacturing.

### Agricultural Production Value

From 1870 to 1910, each Census reports the value of agricultural products at the state level,  $Y_{it}^{ag}$ . We would prefer explicit data on agricultural value added rather than agricultural products. However, in the only year of overlapping values, 1880, the Census numbers match the agricultural income reported by Easterlin in *Trends in the American Economy in the Nineteenth Century*. To determine the state values of agricultural production for 1850, and 1860, we estimate the relationship of the production value of agricultural products sold within a state on the total value of farmland and buildings and agricultural labor force.

Agricultural labor force is reported in the Census in 1840, 1850, and 1870 through 2000. While the census does report a measure of the agricultural labor force in 1850, its usefulness is diminished because it does not include slave labor.<sup>26</sup> To estimate the total agricultural labor force for 1850 and 1860, we use the agricultural labor force reported in 1840, which includes slaves, and in 1870, which includes freed slaves, to construct the portion of the state labor force engaged in agricultural produc-

---

<sup>26</sup>The 1860 census reports data hundreds of detailed occupations, but we do not attempt to map these occupations into the broader agricultural labor force.

tion,  $\text{fraction}_{it}^{ag}$ . In non-slave holding regions, where the omission of slave labor is not problematic, we calculate  $\text{fraction}_{it}^{ag}$  in 1850 using the Census data.<sup>27</sup> We then linearly interpolate  $\text{fraction}_{it}^{ag}$  between 1840 and 1870 (between 1850 and 1870 for slave-holding regions and New England). We complete our measure of agricultural labor force in these intervening years by multiplying  $\text{fraction}_{it}^{ag}$  by the total labor force in each state.<sup>28</sup>

For the 1850 and 1860 values of agricultural products, we estimate the relationship in 1870 and 1880. For 1920, we estimate the relationship in 1910 and 1930.<sup>29</sup> The Census reports the production value of agricultural products and data on total farmland value comes from HSUS. With our measures of agricultural capital,  $\text{farmvalue}_{it}$ , and labor,  $\text{aglabor}_{it}$ , we estimate the value of products produced in 1850, 1860, and 1920 by regressing the following:

$$\ln(Y_{it}^{ag}) = \beta_1 \ln(\text{farmvalue}_{it}) + \beta_2 \ln(\text{aglabor}_{it}) + \beta_3 Z + \beta_4 \text{year}_t \quad (93)$$

where  $Z$  is the vector of region dummies and  $\text{year}_t$  is a time trend. We then take the exponential of the predicted value,  $\widehat{Y}_{it}^{ag}$ , to estimate state level agricultural production value for 1850, 1860, and 1920.

## Manufacturing Value Added

The value added by manufacturers at the state level,  $Y_{it}^{manu}$ , is calculated by subtracting the value of materials used from the value of products sold reported in the Census from 1850 through 1920. Because the 1840 Census does not report the value added by manufacturing, we use the relationship between value added and the manufacturing labor force from 1850 through 1860 to determine value added in 1840. We regress the natural log of value added in the manufacturing sector,  $\text{mvalue}_{it}$ , on the natural log of the manufacturing labor force,  $\text{mlabor}_{it}$ , interacted with regions as well as individual census region effects,  $Z$ :<sup>30</sup>

---

<sup>27</sup>These regions are the Middle Atlantic, Mountain, Pacific, East North Central, and West North Central regions. We do not include the New England region because data in 1850 appear unreliable.

<sup>28</sup>No data on agricultural labor force is reported for Kansas, Nebraska, Texas, and Washington in 1840, therefore, we are unable to calculate the fraction of the labor force in agriculture using the methodology described above. For 1860, we proxy the agricultural labor force for these states by the number of persons listing their occupation as farmers.

<sup>29</sup>Additionally, data on agricultural products is not available in Arizona and New Mexico in 1890. We again regress using Eq. 94 and use data from 1880 and 1900 to estimate values for these two states.

<sup>30</sup>Data on manufacturing labor are not available in 1890 and 1910. We calculate the fraction of the labor force engaged in manufacturing,  $\text{fraction}_{it}^{min}$  in 1880, 1900, and 1920. We linearly interpolate the value of  $\text{fraction}_{it}^{min}$  in 1890 and 1910, and multiply the result by the total labor force.

$$\ln(mvalue_{it}) = \beta_1 Z + \beta_2 (Z \ln(mlabor_{it})) + \beta_3 year_t \quad (94)$$

Taking the exponential of the predicted  $\ln(\widehat{mvalue}_{it})$  generates the 1840 estimate of value added by manufacturing.

## Mining Value Added

The output of precious metals is an important component of state income in the Pacific and Mountain region, particularly so in the early portion of our data set. As will be discussed in the following section, our income calculations allow for a component of income not captured by agriculture and mining. However, our methodology implicitly assumes that this component is relatively stable over time. Given the nature of gold and silver discoveries and subsequent rushes, we find this assumption unsatisfactory for these regions. As a result, we have collected data on precious metals mining output for the Mountain and Pacific regions.

Value added in the precious metals mining sector of the economy is calculated by subtracting the value of materials from the value of mining products,  $product\_value_{it}$ , where available. A measure of mining products is available at the state level from the 1890 Census Report on Mineral Industries in the United States for 1870, 1880, and 1890.<sup>31</sup> A measure of materials used and labor is also available. This allows a measure of mining value added in 1890,  $Y_{i,1890}^{mn}$ , to be calculated.

$$Y_{i,1890}^{mn} = product\_value_{it} - materials_{it} \quad (95)$$

We next calculate per worker value added in 1890:

$$y_{i,1890}^{mn} = \frac{Y_{i,1890}^{mn}}{L_{i,1890}^{mn}} \quad (96)$$

and fraction of output this is value added,  $fracY_{i,1890}$ :

$$fracY_{i,1890} = \frac{Y_{i,1890}}{product\_value_{i,1890}} \quad (97)$$

The 1870 Census report, *The Statistics of Mining*, gives data on employment, materials, and output of precious metals in 1870, but appears to be only a partial sample of all mining establishments. We do not use the measures of total products, value added and employment, but maintain measures

---

<sup>31</sup>Data is not readily available from this source for 1890. Instead, we use the values in 1889

of *per worker* products, value added, and employment.<sup>32</sup> Thus, we calculate  $y_{i,1870}^{mn}$  and  $fracY_{i,1870}$  and then use these values with the 1890 values to interpolate to obtain  $y_{i,1880}^{mn}$  and  $fracY_{i,1880}$ . Prior to 1870, data is not as detailed. We assume that products per worker for each state in 1850 and 1860 is equal to it's value in 1870.<sup>33</sup> Thus:

$$y_{i,1850}^{mn} = y_{i,1860}^{mn} = y_{i,1870}^{mn} \quad (98)$$

We do the same for the fraction of products that is value added.

$$fracY_{i,1850}^{mn} = fracY_{i,1860}^{mn} = fracY_{i,1870}^{mn} \quad (99)$$

We next turn our attention to employment in precious metals mining. Direct measures of precious metals mining employment are available in 1840, and 1890 (and in 1870 we have a sample), as are measures of non-precious metal mining employment. This overlapping data will be exploited below. Data on precious metals employment data do not exist directly in 1850, 1860, and 1880, yet measures of total employment in mining (precious and non-precious) are available in these years.

Let employment in precious metals mining be  $L_{it}^{prec}$ , and employment in non-precious metals mining,  $L_{it}^{nonprec}$ . In 1840, 1870, and 1890 we calculate:

$$fracL_{it}^{prec} = \frac{L_{it}^{prec}}{(L_{it}^{prec} + L_{it}^{nonprec})} \quad (100)$$

For states in which we have no data prior to 1870, we assume that  $fracL_{it}^{prec}$  in 1850 and 1860 are identical to the 1870 values in each state. We also interpolate between 1870 and 1890 to acquire 1880 values. Thus:

$$fracL_{i,1850}^{prec} = fracL_{i,1860}^{prec} = fracL_{i,1870}^{prec} \quad (101)$$

Next, we calculate labor in the precious metal sector,  $L_{it}^{prec}$ , in 1850, 1860, and 1880 as,

$$L_{it}^{prec} = fracL_{it}^{prec} \left( fracL_{it}^{prec \& nonprec} \right) \quad (102)$$

And to correct for the fact that  $L_{it}^{prec}$  in 1870 is a sample, we geometrically interpolate between the value of  $L_{it}^{prec}$  in 1860 and 1880.

---

<sup>32</sup>In addition, we maintain the fraction of all mining labor that is engaged in precious metals mining. See below.

<sup>33</sup>There is only one state, California, for which we have data in 1850. We make a separate adjustment for this state below.

Finally, we can calculate our measure of  $Y_{it}^{mn}$  for 1850, 1860, 1870, and 1880:

$$Y_{it}^{mn} = y_{it}^{mn} L_{it}^{mn} \text{frac} L_{it}^{prec} \quad (103)$$

As a check on the reasonableness of our calculations, we compare the sum of mining output across the states to the national output figures given for 1850 and 1860 in the 1890 Census report. We find we overestimate mining output in 1860. We assume that California has the same share of national mining output in 1860 as it does in 1850. We then renormalize all other states so that the sum is equal to the national total.

### Total State Income

Adding the value of products produced by manufacturers and mines and the estimated income from agricultural production at the state level generates the total state income attributable to manufacturing, mining, and agriculture:

$$Y_{it}^{ag+manu+mn} = Y_{it}^{ag} + Y_{it}^{manu} + Y_{it}^{mn} \quad (104)$$

for  $1840 \leq t \leq 1920$ .<sup>34</sup>

Unfortunately for us, this measure of income is not the total state income, but only the portion of state income resulting from manufacturing, mining, and agriculture. In order to account for the remaining industries in a states' economy, we turn to the total income calculations reported by Easterlin. In *Trends in the American Economy in the Nineteenth Century*, Easterlin calculates the total state income level for 1840 and in *Analyses of Economic Change in Population Redistribution and Economic Growth, United States, 1870-1950*, he reports total state income for 1880, 1900, and 1919-1921(1920). For 1840, 1880, 1900, and 1920, we calculate the difference between our estimated,  $Y_{it}^{ag+manu+mn}$ , and Easterlin's total state income,  $Y_{it}^E$ :

$$Y_{it}^{not} = Y_{it}^E - Y_{it}^{ag+manu+mn} \quad (105)$$

for  $t=1840, 1880, 1900, \text{ and } 1920$ . We then calculate the ratio of income generated outside agriculture, manufacturing, and mining over income produced by agriculture, manufacturing, and mining.<sup>35</sup>

---

<sup>34</sup>We only make our mining adjustments in 1850, 1860, 1870, and 1890 for the Mountain and Pacific regions. We do not adjust mining for states outside of these regions. That is,  $Y_{it}^{mn} = 0$  for all other regions.

<sup>35</sup>We occasionally observe a measure of  $Y_{it}^{not}$  that is less than zero in 1840. For these states, the sum of agricultural,

$$Y_{it}^{notshare} = \frac{Y_{it}^{not}}{Y_{it}^{ag+manu+mn}} \quad (106)$$

For the states with 1840 Easterlin incomes, listed in Table C1, we estimate the ratio of income generated outside agriculture, manufacturing, and mining over income produced by agriculture, manufacturing, and mining for 1850, 1860, 1870, 1890, and 1910 using the following methods:

$$\widehat{Y}_{i,1850}^{notshare} = (Y_{i,1840}^{notshare})^{.75} (Y_{i,1880}^{notshare})^{.25} \quad (107)$$

$$\widehat{Y}_{i,1860}^{notshare} = (Y_{i,1840}^{notshare})^{.5} (Y_{i,1880}^{notshare})^{.5} \quad (108)$$

$$\widehat{Y}_{i,1870}^{notshare} = (Y_{i,1840}^{notshare})^{.25} (Y_{i,1880}^{notshare})^{.75} \quad (109)$$

$$\widehat{Y}_{i,1890}^{notshare} = (Y_{i,1880}^{notshare})^{.5} (Y_{i,1900}^{notshare})^{.5} \quad (110)$$

$$\widehat{Y}_{i,1910}^{notshare} = (Y_{i,1900}^{notshare})^{.5} (Y_{i,1920}^{notshare})^{.5} \quad (111)$$

For the states without 1840 incomes, listed in Table C2, we use the 1880 ratio of income generated outside agriculture, manufacturing, and mining over income produced by agriculture, manufacturing, and mining,  $Y_{i,1880}^{notshare}$ , in order to determine  $Y_{i,t}^{notshare}$ , for t =1850, 1860, 1870. For 1890, and 1910 we use the similar method as above:

$$\widehat{Y}_{i,1850}^{notshare} = (Y_{i,1880}^{notshare}) \quad (112)$$

$$\widehat{Y}_{i,1860}^{notshare} = (Y_{i,1880}^{notshare}) \quad (113)$$

$$\widehat{Y}_{i,1870}^{notshare} = (Y_{i,1880}^{notshare}) \quad (114)$$

$$\widehat{Y}_{i,1890}^{notshare} = (Y_{i,1880}^{notshare})^{.5} (Y_{i,1900}^{notshare})^{.5} \quad (115)$$

$$\widehat{Y}_{i,1910}^{notshare} = (Y_{i,1900}^{notshare})^{.5} (Y_{i,1920}^{notshare})^{.5} \quad (116)$$

Using these ratios we calculate our final total state income,  $\widehat{Y}_{it}^{all}$ , for all non-Easterlin years:

$$\widehat{Y}_{it}^{all} = Y_{it}^{ag+manu+mn} \left[ 1 + \widehat{Y}_{i,t}^{notshare} \right] \quad (117)$$

In order of find our calculated per worker income, we simple take total state income in year and divide it by the states' labor force reported by the census, except 1850 and 1860 where the our labor mining, and manufacturing income exceeds the figure given as total income by Easterlin. We replace the measure of  $Y_{it}^{not}$  with zero. Cases are rare and magnitudes are small.

force figures are adjusted for slaves:

$$y_{it} = \frac{\widehat{Y}_{it}^{all}}{L_{it}} \quad (118)$$

We then put our per worker income measures into real terms by adjusting for both national and regional differences in prices. See Appendix B for more details on price levels.

Table C1: 1840 State Incomes Reported By Easterlin

Alabama	Iowa	Mississippi	Pennsylvania
Arkansas	Kentucky	Missouri	Rhode Island
Connecticut	Louisiana	New Hampshire	South Carolina
Delaware	Maine	New Jersey	Tennessee
Florida	Maryland	New York	Vermont
Georgia	Massachusetts	North Carolina	Virginia
Illinois	Michigan	Ohio	Wisconsin
Indiana			

Table C2: 1840 State Incomes Not Reported By Easterlin

(with first year of agriculture and manufacturing data availability)

State	First Year Calculated	State	First Year Calculated
Arizona	1870	New Mexico	1850
California	1850	Oregon	1850
Colorado	1870	South Dakota	1910
Idaho	1870	Texas	1850
Kansas	1860	Utah	1850
Minnesota	1860	Washington	1860
Montana	1870	West Virginia	1870
Nebraska	1860	Wyoming	1870
Nevada	1870		

## APPENDIX D

Table D1 below presents the correlations of our years of schooling in the labor force with the two separate state human capital measures of Mulligan and Sala-i-Martin (1997,2000).

D1: Correlation of Years of Schooling in the Labor Force with

Mulligan and Sala-i-Martin (1997, 2000)			
1940	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.9326	1	
hc2000	.8996	.9747	1
1950	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.8824	1	
hc2000	.8081	.9321	1
1960	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.7766	1	
hc2000	.7955	.9500	1
1970	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.6455	1	
hc2000	.6727	.8403	1
1980	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.8466	1	
hc2000	.7669	.8792	1
1990	yrs of schooling	hc1997	hc2000
yrs of schooling	1		
hc1997	.8449	1	
hc2000	.7797	.9141	1

Table D2 below details how well we fit the census information using labor force weighted regressions.<sup>36</sup> Overall, the our calculations fit the data extremely well, but this could be due to the trend in education. Thus we present the decade by decade results. If our estimates were exactly in line with the census, we would get a slope coefficient of 1 on years and a 0 intercept.

<sup>36</sup>This seems reasonable as it seems much more important to fit New York or California than to give those states equal weight with states like North and South Dakota.



Table D2: Regressions of Average Years of Schooling from the Census on Estimates

(standard errors)								
variable	ALL	1940	1950	1960	1970	1980	1990	2000
$E$	1.053	1.028	1.147	1.185	1.161	0.948	0.8101	0.855
	(0.009)	(0.002)	(0.002)	(0.007)	(0.007)	(0.003)	(0.003)	(0.003)
constant	-0.565	-0.355	-1.449	-1.833	-1.835	0.738	2.416	1.961
	(0.103)	(0.145)	(0.193)	(0.693)	(0.830)	(0.462)	(0.505)	(0.461)
$N$	355	49	51	51	51	51	51	51
$\bar{R}^2$	.9723	.9127	.9183	.7963	.7889	.8413	.7982	.8425
$prob > F$	.0000	.0013	.0001	.0973	.0034	.0001	.0002	.0000

The final row of the table contains the result of the joint test of this hypothesis. Overall we reject the null hypothesis that our estimated slope coefficient is 1 and our intercept is 0, however for 1960 we cannot reject the null. In all regressions, our fit is quite good, with  $\bar{R}^2$  over .75.

An alternative way to compare our estimates of years of schooling in the labor force with the values of years of schooling by state from the Census is to compare the means and standard deviations weighted and unweighted. Table D3 provides evidence that our estimates are similar, if not identical with the census values.

Table D3: Average Years of Schooling: Census and New Estimates

year	Census mean	Census std. dev.	Estimate mean	Estimate std. dev.	% dev. mean	Census weighted mean	Estimate weighted mean	% dev. weighted mean
1940	8.24	0.89	8.34	0.86	1.2	8.17	8.29	1.3
1950	8.98	0.87	9.07	0.74	1.0	8.95	9.07	1.3
1960	9.85	0.75	9.85	0.58	0.0	9.82	9.83	0.1
1970	10.68	0.64	10.85	0.49	1.6	10.65	10.75	0.9
1980	11.87	0.59	11.65	0.51	-1.9	11.82	11.68	-1.2
1990	12.45	0.42	12.37	0.43	-0.6	12.43	12.36	-0.6
2000	13.14	0.38	13.04	0.41	-0.8	13.08	13.01	-0.5

Table D3 shows that our average years of schooling measure nearly match Census estimates. The largest weighted difference occurs in 1940 and 1950, while the largest unweighted difference occurs in 1980. The smallest difference occurs in 1960 for both measures. From 1940 onward the mean of

our estimates differs from the Census by less than 1.9 percent. One thing evident from Table D3 is the greater amount of dispersion about the mean in our estimates from 1980 to 2000, but smaller dispersion than the Census estimates before 1980.