

# Education Benefits of Universal Primary Education Program: Evidence from Tanzania

Esther DELESALLE \*

October 25, 2016

## Abstract

The purpose of this paper is to determine the impact of education on labor market participation and on household's consumption depending on the sector of activity. To do so, I refer to the UPE program implemented in Tanzania from 1974 to 1978. The exposition to this program varied according to the year of birth and the region of birth. Thus, I exploit these two exogenous variations to instrument education of the household's head. I find that the UPE program decreased inequalities of access to education by increasing education of less educated regions. This raise in education has a positive impact on household's consumption aggregates. The magnitude of this effect varies between sectors of activity and geographical areas : the effect is smaller in agriculture than for non-agricultural self-employment activities and for wage-earning activities, and the effect in urban areas is twice as large as the effect in rural areas.

**Keywords:** Human capital investment, Returns to Education, schooling reforms, Tanzania.

**JEL Codes:** XXX

---

\*University of Cergy-Pontoise, THEMA, 33 Boulevard du Port, F-95011 Cergy-Pontoise and Paris School of Economics, 48 boulevard Jourdan 75014 Paris, France . Email: estherdelesalle@gmail.com

# 1 Introduction

Education is a cornerstone for economic growth and plays a crucial role in labor markets. As a consequence, governments and non-governmental organisations have put education at the top of their agenda. More specifically, several governments of developing countries have implemented policies to universalizing primary education. An extensive body of literature underlines the positive correlation between education and earnings but does not inform about the causality of the relationship. Card [2001] reviews papers that aim to identify the causal impact of education on earnings. To disentangle between the ability effect and the education effect on earnings, these selected papers either instrument education based on characteristics of the schooling system, or they use family background as control or instrument. Among the eleven papers including in the survey, only two of them focus on developing countries, the Duflo [2000]' paper where education is instrumented by a school construction program in Indonesia and the Maluccio [1998]'s paper where education is instrumented by the distance to school in rural Philippines. Both authors restrict their analysis to individuals who earn a wage. Thus, it raises the question of the representativity of these samples. Indeed, wage-earnings individuals are likely to be self-selected and to have specific characteristics. Maluccio [1998] does not deal with this sample selection issue while Duflo [2000] adopts an imputation technique to compute a wage for individuals from the self-employment sector. She finds that returns to education significantly drop with this method. This method is compatible for countries with a developed formal sector but it is less adapted for Sub-Saharan countries that are mainly agriculture-based countries. Indeed, in agriculture, the production is available at the household's level and few individuals are wage-earners.

To answer to this substantial issue, other papers estimate the returns to education at the household' level. [Griliches, 1964] is the first researcher to measure the impact of education of the household's head on agriculture with a joint production function. Lockheed et al. [1980] review papers estimating the impact of education on agricultural production and find very mixed results depending on the country of reference and the specification of education. However, the latter do not take into account the endogeneity of education of the household's head.

The first contribution of this paper is to assess the efficiency of a massive primary education program. In line with this, I evaluate whether the Universalization Primary Education program implemented in Tanzania from 1974 to 1978 ensures the expansion of the education system and contributes to reduce inequality with regards to access to education.

The second contribution of this paper is to estimate the returns to education in developing countries for the entire population when education is considered as an endogenous variable. Since developing countries are often

characterized by the supremacy of the informal sector and the agricultural sector, I use consumption aggregates that are available for all sample households.

Given that education may be endogenous, I instrument education of the household's head by the UPE program that constitutes a natural experiment. In 1974, educational levels were low at the national level (5 years of education) with strong variations between regions. The introduction of the UPE program led to substantial results: 3,3 million of children aged 7 to 13 were enrolled in 1980 compared to 1,2 million in 1974 [Bonini, 2003]. To reduce disparities of access to education, the Tanzanian communist government put a lot of emphasis on deprived areas, which led the latter to experience higher schooling expansion. Therefore, the exposure to the UPE program varied according to two criteria: the age of the individual at the time of the reform and the educational level of the region before the introduction of the program. The UPE program gave rise to an exogenous variation in education and should be a valid instrument for education. Thus, I rely on this instrument to determine the effect of education on consumption. In order to capture the variability of the returns to education, I also distinguish the returns to education for subgroups: rural areas, urban areas, the agricultural sector, the self-employment sector and the formal sector.

The third contribution of this paper is to address the effect of education on the probability of working in a given sector of activity. In order to do so, I adopt the same identification strategy and I instrument education by exploiting the nature of the UPE program.

The main finding of this paper is that UPE program reduces inequalities of access to education and that education augmented consumption aggregates from 7.4 percent to 10.5 percent. I notice, however, that these returns to education seem to be lower in rural areas compared to urban areas. Similarly, returns to education estimates are lower for the agricultural sector than for the non-agricultural self-employment sector<sup>1</sup> and for the wage-earner sector.

The remainder of the paper is organized as follows: section 2 provides a broad picture of the evolution of education in Tanzania and describes the data and the main variables of the analysis. Section 3 introduces the identification strategy; section 4 presents the effect of the UPE program on education; section 5 and section 6 respectively focus on the effect of education on consumption and on the labor market's participation. Section 7 deals with sample selection issues while section 8 concludes.

---

<sup>1</sup>With and without employee

## 2 The program

### 2.1 Historical background and the UPE program.

After the end of colonization, access to education in Tanzania varied a lot from one region to another [Kinyanjui et al., 1980]. These spatial disparities based on ecological endowment were exacerbated by colonial activities and transport networks.<sup>2</sup> The arrival in power of the Prime Minister Nyerere in 1964 marked a radical political and economic change. The policy of Education for Self Reliance (ESR) was decreed in 1967 at the Arusha conference. Education became the mainstay of the Tanzanian socialist economy that should insure economic growth and the Tanzanian autonomy. Afterwards, the government started to allocate more funding to education.

In 1974, the government reaffirmed this will and committed itself to reach the Universal Primary Education (UPE) by 1978. The aim was to improve the equity of access to education and to teach agricultural skills that insure efficiency and self-reliance of rural communities. During this reform, the education administration was at the regional level but wards and villages were responsible for the schools' organization. To achieve UPE, the government made series of changes to develop the schooling infrastructures in order to welcome all children aged 7 to 13.

The forced villagization eased the increase of education attainments. Households living in remote areas were forced to move in community villages called ujamaa. Most of the time, the distance to the prior dwelling was lower than five kilometers. The aim was to gather individuals to provide them all social services, including schooling. From 1974 to 1977, more than 10 millions of people were moved and 2650 ujamaa were built [Martin, 1988]. Prime Minister Nyerere [1987] considered that gathering the rural population was necessary to develop education, to reduce inequalities and to improve agricultural production. Besides, the Tanzanian government massively invested in primary education and concentrated its efforts on deprived areas. Local resources were mobilized for classrooms and a large number of new schools were built. In 1978, expenditures on primary education was three times the amount dedicated to secondary education [Bonini, 2003]. Then, the UPE program combined with villagization largely contributed to reduce distance to school.

Simultaneously, teachers' recruitment and teachers' training were restructured to deal with the growing number of pupils. The Primary Education Reform Kwamsisi Project in 1970 was designed to anticipate this. It improved the skills of 10,000 teachers. Despite this program, there was still a shortage of primary school teachers that may have affected the quality of education, especially in the beginning of the UPE plan [Sabates et al., 2011].

---

<sup>2</sup>The most privileged zones were the Arusha-Kilimanjaro-Tanga corridor, the Coast Morogoro-Kigoma and the Mwanza-Shinyanga corridor [Maro and Mlay, 1979].

The government made additional adjustments to improve schools' attractiveness. Tuition fees were suppressed, primary education became mandatory and the language of instruction became the Swahili, the mother tongue of most pupils.

The aim of primary school was also to train pupils to become potential farmers. To fulfill this goal, the examination in the middle of the primary cycle was removed, the starting age was postponed from 5 to 7 years old and agriculture classes were introduced in the curriculum. As a result, pupils leaving the primary schools would have been old enough and would have acquired the abilities to work in the fields. To encourage people to start working after primary school, access of the secondary cycle was drastically limited by regional quotas [Martin, 1988].

Results of this UPE plan were considerable: 90 percent of children aged 7 to 13 were enrolled in 1978, access to primary education was improved and disparities among regions were reduced [Bonini, 2003].

## **2.2 Data**

### **2.2.1 Data set**

The Data of this paper come from two sources: the 10 percent IPUMS sample from the 2002 Tanzanian Census and the LSMS-ISA (LSMS-Integrated Surveys on Agriculture) panel data collected in 2008-2009, 2010-2011 and 2012-2013<sup>3</sup>. Both surveys were designed to be nationally representative. The IPUMS data contain 809,699 households and the LSMS-ISA data include 3265 households in 2008, 3924 households in 2010 and 5010 households in 2012.<sup>4</sup> The sample of my analysis is composed by households' head born between 1945 and 1954 and between 1961 and 1978. The 2002 census subsample includes 500,519 households and the LSMS-ISA data contains 3,706 households.

### **2.2.2 Main variables**

Education is one of the main interest variable of the analysis. I consider two variables, the number of years of education (between 0 and 15 years of education) and a dummy variable indicating whether the individual has completed primary education.

To measure the returns to education, several outcomes can be used. Living standards can either be measured with income or consumption. In this analysis, consumption presents several advantages. First, consumption is less subject to variations than income in rural agriculture. This variable is

---

<sup>3</sup>From October 2008 to December 2009 for the first wave, from October 2010 to December 2011 for the second wave, and from October 2013 to December 2013 for the third wave.

<sup>4</sup>The number of households is increasing over the three waves due to the split-off.

less affected by shocks and should be more representative of household's well-being [Deaton and Zaidi, 2002]. Second, income is not measured in the same way between non-agricultural self-employed activities and agriculture, the two prevailing activities in Tanzania.<sup>5</sup> It puts into question the reliability of the comparison. On the contrary, consumption aggregates is similarly defined for all households whatever their sector of activity. Third and last, consumption can be computed for all households. Thus, it is a way of avoiding selection and imputation issues.

[Deaton and Zaidi, 2002] propose a detailed guideline to construct consumption aggregates from household survey data. I adopt their method that is particularly suitable for LSMS data and I define two per capita consumption variables composed by four sub-aggregates, food items, non food items (education expenses, health insurance, etc), housing consumption (rent expenditure and consumption of utilities) and consumer durables. Both of these consumption aggregates are adjusted by a price index to take into account variations in prices faced by households. These two indexes differ by the fact that one of them is divided by an equivalence scale<sup>6</sup>.

However, accurate consumption data are not available in all household's survey and require large resources. An alternative approach is to build asset indexes based on assets, access to utilities and housing characteristics [Vyas and Kumaranayake, 2006]. Construction of this index is obtained by factor analysis<sup>7</sup>. A large number of papers ([McKenzie, 2003], [Vyas and Kumaranayake, 2006], etc) find that these indexes correctly measure inequalities between households and are good proxies for long-term wealth. The additional advantage of these indexes is that it limits measurement errors [Sahn and Stifel, 2003]. In the IPUMS data, some dwelling characteristics and utilities are reported<sup>8</sup>. By relying on these dimensions and by adopting a factor analysis, I construct a wealth index that should be a good proxy of household's welfare.

### 3 Identification strategy

The UPE program was applied during a limited time frame and was mainly targeting regions with poor access to education, hence the exposure to the program is captured by two dimensions: region of birth and year of birth.

Since individuals could have moved from a region to another during the

---

<sup>5</sup>Self-employment income is rarely a wage and agricultural income can be captured through the production.

<sup>6</sup>the equivalence scale is made from the household's size when every adult represents one unit and each child represents 0.3 unit to allow for the fact that children consume less than adults

<sup>7</sup>The aim is to convert a set of variables into a smaller number of indices.

<sup>8</sup>whether the household has drinking water, electricity, a phone, a flush toilet, a solid roof, solid walls, etc.

UPE program, regions of residence are endogenous to the program. Thus, I refer to the region of birth that has been determined prior to the program.

The Tanzanian primary cycle starts at 7 and ends at 13 years old. It implies that children older than 13 in 1974 have not been exposed to the reform. It is worth noting that pilot programs started by 1968, just after the Arusha conference: some regions benefited from financial supports and the villagization procedure started sooner in some areas.<sup>9</sup> Besides, overage children could have taken advantage of the reform after 13 if they were enrolled in delay or if they had repeated classes. As a result, individuals who were younger than 13 between 1968 and 1974 were likely to be partially treated. To avoid contamination issues, I define the pre-treatment group T0 such as it only contains household heads not affected by the UPE reform. It is composed by individuals born between 1945 and 1954 that were too old to go back to primary education at the time of the reform. Several age-cohorts have been affected by the UPE plan with a variable intensity. According to the age of the household's head, I distinguish three treated groups (see Table 1). T1 contains households with heads aged 8 to 13 in 1974. The second treatment group T2 includes households with heads aged 3 to 7 in 1974. All individuals from T1 and T2 have been concerned by the reform before it was stopped in 1978: individuals from T1 were treated at the beginning of the reform while individuals from T2 were treated at the end of the reform. T3 gathers individuals aged 0 to 6 in 1978. This group can be used to test whether the effect of the UPE program was persistent over time.<sup>10</sup>

Table 1: Age Cohorts

Age cohorts	Year of birth	Age in 1974	Potential education level during the UPE plan	Obs. IPUMS	Obs. LSMS
T0	1945-1954	20- 29	postsecondary and over	111,818	1,706
T1	1961-1966	8 -13	primary-secondary	113,063	1,554
T2	1967-1971	3-7	no education- primary	103,406	1,407
T3	1972-1978	. - 2 (0-6 in 1978)	no education	172,232	2,154

Figures 5 in Appendix shows the education distribution among these age-cohorts. One notices that the percentage of individuals with no education drastically decreased from T0 to T1 (it moved from 45 to 24 percent): not only more people enrolled to primary, but more people completed primary education. The percentage of primary completion kept growing from T1 to T2 and increased by 10 percent. However, the percentage of people with no education remains stable. Although there is a substantial increase in education over time, a non-negligible share of the population did not enrol.

<sup>9</sup>The Tanga region in 1968, the west lake region in 1968, Dodoma in 1969, the Iringa region in 1971 (Martin, 1988).

<sup>10</sup>schools have not been destroyed and a significant number of new teachers were employed.

Figure 1 shows the relationship between the regional education attainment in T0 and the increase in education between T0 and T1. Each dot corresponds to a region of birth. One notes that the UPE program has been more intense in regions with low schooling enrollment at T0. Thus, it should have helped reducing disparities in access to education.

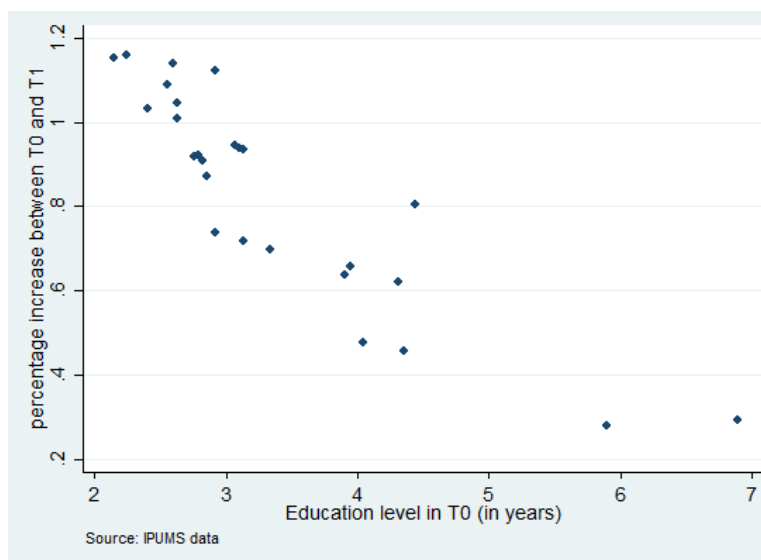


Figure 1: Evolution of education attainment by region from T0 to T1 according to the education level in T0.

However, it is important to underline that this relation is not necessarily caused by the UPE program. In the case of convergence, less developed regions could have had a higher education increase in order to catch up the more advanced regions. If this were to be true, this phenomenon would be observed before and after the introduction of the UPE program. To test this assumption, I compare the education trend from a baseline period T<sub>b</sub> to T<sub>0</sub> and from T<sub>0</sub> to T<sub>1</sub>. T<sub>b</sub> is delineated just before T<sub>0</sub> and includes individuals born from 1935 to 1944. I gather regions in two groups: *region +* had a high education attainment in 1958 and *region -* had a low education attainment in 1958<sup>11</sup>. Table 2 shows that education increases overtime for both regions. From the baseline period to T<sub>0</sub>, one notices that this increase is higher in *region +*: education attainment is raised by 3.8 years whereas it is only raised by 2.5 years in *region -*. Then, instead of being reduced, inequities of access to education are worsen by 1.29 year. This result can be explained by the fact that *region +* are likely to have on average larger wealth and to be more efficient at improving education. From T<sub>0</sub> to T<sub>1</sub>, this trend is reversed. A higher increase is observed in *region -*, which contributes

<sup>11</sup>It corresponds to the year when the first age-cohort started primary education.



to significantly reduce the education gap by 0.386 year. These differences in differences suggest that the introduction of the UPE program allowed a greater education expansion in remote areas.

Table 2: Evolution of education attainment by period and region groups

Age-cohort	Region -	Region +	Difference (region+ - region-)
Tb	1.536 (.0119)	2.806 (.0227)	1.270 (.0270)
T0	2.571 (.0128)	4.249 (.0221)	1.677 (.0255)
T1	4.746 (.0128)	6.265 (.0186)	1.518 (.0226)
T2	4.524 (.0119)	6.675 (.0167)	1.150 (.0205)
T3	5.392 (.008)	6.529 (.0119)	1.137 (.0144)
Difference (T0-Tbaseline)	1.035 (.0175)	1.443 (.0316)	.408 (.0362)
Difference (T1-T0)	2.175 (.0181)	2.016 (.0288)	-.158 (.0340)
Difference (T2-T0)	2.952 (.0175)	2.426 (.0276)	-.527 (.0328)
Difference (T3-T0)	2.821 (.0152)	2.280 (.0251)	-.541 (.0293)

Note: Source: IPUMS data, 2002. Standards errors are in parenthesis. Region + represents regions with education in 1958 higher than 3 years of education, the average primary education level by this time. Region - represents regions with education in 1968 lower than 3 years of education, the average primary education level by this time. Tb: individuals born between 1935 and 1945.

## 4 The impact of the UPE program on the education expansion

Exposure to the UPE program varies according to the region of birth and the year of birth. I use this difference in difference to measure the impact of the UPE program on education:

$$S_{ijt} = \alpha + \beta_j + \beta_t + \theta T * S_{j,1967} + \delta T * X_{j,1967} + \epsilon_{ijt} \quad (1)$$

$\beta_j$  and  $\beta_t$  are region-of-birth fixed effects and birth-cohort fixed effects to account for permanent differences across regions and over time.  $T$  is a dummy taking value 0 for people belonging to T0 and 1 for people belonging to T1, T2 or T3.  $S_{j,1967}$  denotes the average education level in region  $j$  in 1967. To predict the intensity of the UPE plan, I refer to  $S_{j,1967}$ , the average education level in region  $j$  in 1967, before the introduction of the UPE program and of pilot programs.

Given that the primary school ends at 13 years old, I compute  $S_{j,1967}$  with the education level of individuals born in region  $j$  in 1954 to predict the average education level in 1967.  $X_{j,1954}$  is a set of region' characteristics. It includes the population aged 7 to 13 in 1967 and the percentage of rural areas in 1967.  $S_{j,1967}$   $\theta$  captures the effect on education of the expansion of the education level between T0 and T1, T2 or T3 due to the UPE program : when  $S_{j,1945}$  increases, the UPE program's intensity should decrease and the expansion between T0 and treatment groups is expected to be smaller.

Table 3 reports the results. One additional year of initial education level decreased the education expansion by 0.369 years between T0 and T1, by 0.504 years between T0 and T2 and by 0.511 between T0 and T3. F-test values are large and significant for all estimations. This result is consistent with the idea that the UPE program mostly targets regions with low initial education level. Comparison of columns (1-3) with columns (4-6) informs that the program intensity keeps growing from T1 to T2. On the contrary, it seems to stabilize from T2 to T3 (the post-treatment group). Comparison of columns (2-3), (5-6) and (8-9) allows to identify the effect in rural and urban areas<sup>12</sup>. It suggests that the effect of the program is larger in urban areas (it is 0.12 years higher for T1, 0.10 for T2 and 0.14 for T3). Logistical reasons can explain this gap. First, it is probably easier to enforce people to go to school in urban areas than in rural areas. Second, urban areas probably had more infrastructures to mobilize for schools. Thus, the UPE program may have been implemented more efficiently in urban areas.

The second line of table 3 clarifies whether the UPE plan had fully reached its goal by convincing people not only to enrol to school but also to complete primary cycle. It is not obvious since primary lasts seven years

---

<sup>12</sup>The identification of rural and urban areas is made on actual location and not on the place of birth. It could generate a selection bias that would be tackled in section 7.

and that the UPE plan was implemented during four years. Columns (1), (4) and (7) indicates that, when rural and urban areas are taken together, one additional year of  $S_{j,45}$  significantly decreased the primary completion by 3.1 percent for T1, 4.6 for T2 and 4.1 for T3. However, this effect is not significant in rural areas.

Table 3: Effect of the program on the education level: coefficients of  $T * S_{j,58}$ .

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	All	T1 Rural	Urban	All	T2 Rural	Urban	All	T3 Rural	Urban
years of education	- 0.369*** (0.0615)	-0.280*** (0.0675)	-0.400*** (0.0491)	- 0.504*** (0.0646)	- 0.392*** (0.0698)	-0.496*** (0.0470)	-0.511*** (0.0814)	-0.386*** (0.0886)	- 0.533*** (0.0466)
R-squared	0.317	0.242	0.275	0.338	0.276	0.274	0.297	0.232	0.230
F-test	36.02	17.25	66.41	60.98	31.56	111.4	39.35	18.97	131.1
Primary completion	-0.0309** (0.0148)	- 0.00667 (0.0154)	-0.0485*** (0.00872)	-0.0469*** (0.0161)	-0.0199 (0.0172)	- 0.0566*** (0.00733)	-0.0412** (0.0168)	-0.0197 (0.0184)	-0.0457*** (0.00886)
R-squared	0.289	0.257	0.226	0.328	0.305	0.246	0.283	0.245	0.218
F-test	4.390	0.188	30.95	8.474	1.333	59.52	6.057	1.141	26.65
Observations	204,505	122,169	82,336	195,316	114,682	80,634	254,357	145,404	108,953

Note: Source: the IPUMS data. Standard errors are clustered at the region of birth level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity

I also estimate a more detailed regression that determines the effect of the UPE program according to the time exposure to the program:

$$S_{ijt} = \alpha + \beta_j + \beta_t + \sum_{t=1945}^{1954} \gamma_t t * S_{j1945} + \sum_{t=1961}^{1978} \gamma_t * S_{j1945} + \delta t * X_{j1945} + \epsilon_{ijt} \quad (2)$$

In this equation,  $\gamma_t$  indicates age-cohort coefficients. It measures the effect of the reform by age-cohort. The difference between  $\gamma_{t+1}$  and  $\gamma_t$  represents the evolution of education between  $t$  and  $t+1$  due to the initial education level. This equation should complete information of table 2 by checking that results are not driven by a pre-trend during T0. For the pre-treatment group,  $S_{j,45}$  should have no impact on education expansion and  $\gamma_t$  values should be stable and close to 0. On the contrary, if low educated regions benefit more from an education increase,  $\gamma_t$  values should be decreasing for treated groups.

This is precisely what is shown in graph 2. Each coefficient in this table corresponds to the  $\gamma_t$  coefficients of equation 2 (the reference year is the 1945 cohort). Even if  $\gamma_t$  coefficients slightly evolve during T0, they stay close to 0 ( $\gamma_t$  equals 0,02 in 1954). From 1961 to 1978,  $\gamma_t$  coefficients steadily decrease and lose 0.49 point. The most treated cohort is the 1968-cohort where gamma reaches -0.53. This result seems logical: more stricken individuals

by the reform are individuals who were 6 at the time of the reform in 1974. Younger cohorts are still exposed but the intensity does not increase. This graph confirms that the identification strategy is reasonable: the trend was not present before the program and the UPE plan had a significant impact on education for treated cohorts (all coefficients are significant). Thus, if no regional time-varying characteristics correlated to the program's intensity is omitted, these difference-in-difference results should correctly estimate the impact of the UPE plan.

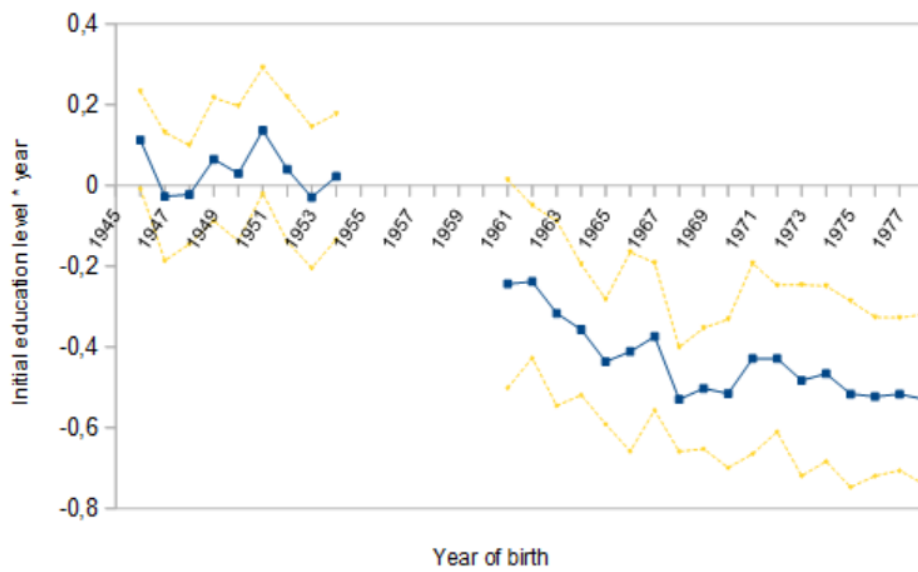


Figure 2:  $\gamma_t$  coefficients of the interaction between age-cohorts and education level by region in 1945. (Source: IPUMS data.)

## 5 The returns to education

In this section, I measure the returns to education by looking at the effect of education  $S_{ijt}$  of the household's  $i$  head born in region  $j$  born at year  $t$  on consumption  $C_{ijt}$ . I confront the results from three consumption indexes, one from the IPUMS data and two from the LSMS data. The main equation is :

$$\text{Log}(C_{ijt}) = \alpha + \beta_j + \beta_t + \Theta * S_{ijt} + \delta T * X_j + \epsilon_{ijt} \quad (3)$$

where  $\beta_j$  and  $\beta_t$  are respectively region-of-birth and year-of-birth fixed effects. Regional controls  $X_j$  are also included . For the sake of comparison of my results, this equation is estimated for the three treatment groups (T1, T2 and T3).

### 5.1 OLS estimations at the individual level

I first ignore the potential endogeneity of education that occurs with OLS regressions. Table 4 gathers results. For the three consumption variables, returns to education do not statically differ between T1, T2 and T3 samples . Regarding the LSMS consumption aggregates, returns to education are about 5.5 percent in rural areas and between 7.5 and 8.5 percent in urban areas. IPUMS estimations are lower: they are about 2.5 percent in rural areas and about 5.8 percent in urban areas. This gap may come from the definition of these variables. The LSMS variables are consumption aggregates whereas the IPUMS variable is an asset index that does not include short-run consumption. Thus, interpretation of these coefficients slightly differs.

These figures are consistent with the literature but they cannot be interpreted as the causal impact of education. If unobserved individual characteristics  $\epsilon_{ijt}$  are correlated with education  $S_{ijt}$ ,  $\Theta$  is biased.

Table 4: OLS estimations of the returns to education

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	All	T1 Rural	Urban	All	T2 Rural	Urban	All	T3 Rural	Urban
IPUMS- Wealth index									
	0.0432*** (0.00254)	0.0248*** (0.00284)	0.0578*** (0.00290)	0.0446*** (0.00245)	0.0263*** (0.00295)	0.0578*** (0.00283)	0.0432*** (0.00227)	0.0253*** (0.00260)	0.0580*** (0.00282)
R2	0.473	0.271	0.391	0.476	0.272	0.381	0.470	0.257	0.345
Obs.	202,196	120,775	81,421	193,221	113,431	79,790	251,674	143,733	107,941
LSMS- Deaton and Zaidi Consumption aggregate (2002)									
	0.0695*** (0.00450)	0.0569*** (0.00552)	0.0756*** (0.00732)	0.0702*** (0.00494)	0.0548*** (0.00628)	0.0814*** (0.00695)	0.0762*** (0.00537)	0.0548*** (0.00536)	0.0847*** (0.00588)
R2	0.484	0.438	0.555	0.446	0.415	0.546	0.473	0.395	0.542
Obs.	2,914	1,499	666	2,779	1,403	636	3,481	1,639	902
LSMS- Deaton and Zaidi Weighted Consumption aggregate (2002)									
	0.0695*** (0.00450)	0.0540*** (0.00545)	0.0774*** (0.00795)	0.0692*** (0.00527)	0.0556*** (0.00568)	0.0797*** (0.00762)	0.0727*** (0.00550)	0.0543*** (0.00502)	0.0802*** (0.00524)
R2	0.484	0.343	0.459	0.429	0.341	0.427	0.461	0.357	0.415
Obs.	2,914	1,499	666	2,779	1,403	636	3,481	1,639	902

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

## 5.2 OLS estimations at the district level

Under the assumption that unobserved individual characteristics varying over time ( $\epsilon_{ijt}$ ) are normally distributed among districts, one can estimate the following district-level equation:

$$\text{Log}(C_{dt}) = \alpha + \beta_d + \beta_t + \Theta * S_{dt} + \epsilon_{dt} \quad (4)$$

$C_{jt}$  is the average consumption of individuals born in district d at time t and  $S_{dt}$  is the related average education level in district d at time t. To obtain unbiased estimates of equation (4), another assumption has to be fulfilled: no district unobserved-variable characteristics that have influenced prior education choices should interfere with the district consumption level observed today. This condition is not respected if there are persistent variable shocks. Results of OLS district-level regressions are presented in Table 5. Coefficients correspond to the effect of one additional year of education in the district on the district consumption level. One can highlight that IPUMS estimates get closer to LSMS estimates: it is about 4 percent in rural areas and estimate range is from 6 to 13 percent in urban areas. If one of the two assumptions does not hold, OLS district-level results are biased. Thus, I confront these results with 2SLS where I instrument education by relying the UPE program.

Table 5: OLS estimations of the returns to education at the district level

	(1)	(2)	(3)	(4)	(5)	(6)
	T1		T2		T3	
	Rural	Urban	Rural	Urban	Rural	Urban
IPUMS- Income index						
	0.0411*	0.0755***	0.0450**	0.135***	0.0409*	0.0615***
	(0.0219)	(0.0133)	(0.0216)	(0.0213)	(0.0228)	(0.0134)
R-squared	0.898	0.897	0.897	0.886	0.897	0.890
Observations	1,990	2,038	1,868	1,915	2,117	2,167
LSMS- Deaton and Zaidi Consumption index						
	0.0687***	0.113***	0.0663***	0.105***	0.0748***	0.118***
	(0.00833)	(0.0132)	(0.00814)	(0.0146)	(0.00757)	(0.0119)
R-squared	0.323	0.383	0.284	0.324	0.296	0.356
Observations	746	402	694	403	817	521
LSMS- Deaton and Zaidi Weighted Consumption index						
	0.0566***	0.0866***	0.0582***	0.0869***	0.0663***	0.0921***
	(0.00753)	(0.0103)	(0.00775)	(0.0126)	(0.00740)	(0.00976)
R-squared	0.314	0.370	0.299	0.331	0.344	0.356
Observations	746	402	694	403	817	521
Cohort FE	YES	YES	YES	YES	YES	YES
District FE	YES	YES	YES	YES	YES	YES

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

### 5.3 2-SLS estimations

The identification strategy is valid if the following assumptions are respected: i) education and consumption by region would not have evolved in the same way in the absence of the UPE program (section 4 provides evidence that it is true), ii) the UPE program does not affect consumption except through education, iii) there is no other regional variable that influences consumption and induces the same time trend reversal than the UPE program. Hence, I use the interaction between the treatment variable T and the initial regional education  $S_{j45}$  to instrument education. Equation (3) that estimates the impact of the UPE program on education becomes the first stage of IV estimates. As shown in Table 3, the UPE program has a high explanatory power (the F-test of the over-identifying restrictions is larger than 10 for all samples and areas). Instead of using a dummy treatment variable as instrument, I can distinguish the effect of the treatment by cohort. I refer to equation (2) but I impose that each  $\gamma_{jt}$  takes value 0 for T0:

$$S_{ijt} = \alpha + \beta_j + \beta_t + \sum_{t=1961}^{1978} \gamma_t * S_{j1945} + \delta_t * X_j + \epsilon_{ijt} \quad (5)$$

By comparing with T0 that is used as a control group,  $\gamma_t$  identifies the UPE program's effect by age-cohort. Table 6 reports 2-SLS estimates for the IPUMS data. I control for the population aged 7 to 13 in 1958 and

the percentage of people living in rural areas in 1958 interacted with time dummies. The second line presents results when the instrument is  $S_{j58} * T$  and the third line presents results when the instrument used is  $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$ . The two 2SLS estimates are similar and show higher returns to education compared to OLS estimates : For the T1 sample when rural and urban areas are taken together, 2SLS estimates are between 1.8 times and 2 times higher than the OLS estimate. It is worth noticing that the gap between rural and urban areas is lower with 2-SLS estimates: returns to education are still higher in urban areas than rural areas but the ratio moves from 2.33 to 1.53 (with  $S_{j58} * T$  instrument) and to 1.49 (with  $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$  instrument). With regard to the ability bias, this result is counter-intuitive. However, as Card [2001] suggests, this common finding could be explained if the instrument used mostly affects low-education subgroups which have higher marginal returns to education. In that case, this effect compensates the ability bias and explains why 2SLS estimate often lead to higher returns to education.

Table 6: Estimations of the education effect on the wealth index (Income index constructed from the IPUMS data)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	All	T1 Rural	Urban	All	T2 Rural	Urban	All	T3 Rural	Urban
OLS	0.0432*** (0.00254)	0.0248*** (0.00284)	0.0578*** (0.00290)	0.0446*** (0.00245)	0.0263*** (0.00295)	0.0578*** (0.00283)	0.0432*** (0.00227)	0.0253*** (0.00260)	0.0580*** (0.00282)
R2	0.473	0.271	0.391	0.476	0.272	0.381	0.470	0.257	0.345
IV: $S_{j58} * T$	0.0861*** (0.0258)	0.0669* (0.0380)	0.0997*** (0.0111)	0.0741*** (0.0197)	0.0487** (0.0218)	0.0882*** (0.0138)	0.105*** (0.0188)	0.105*** (0.0267)	0.0986*** (0.0142)
R2	0.334	0.084	0.241	0.349	0.103	0.259	0.320	-0.006	0.215
F-test	35.69	17.15	65.43	61.10	31.71	110.2	39.31	19.07	131.3
IV: $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$	0.0780*** (0.0268)	0.0614* (0.0373)	0.0942*** (0.0140)	0.0723*** (0.0212)	0.0415* (0.0243)	0.0851*** (0.0202)	0.106*** (0.0190)	0.102*** (0.0258)	0.0969*** (0.0139)
R2	0.342	0.091	0.254	0.350	0.105	0.197	0.322	0.006	0.218
F-test	19.98	10.74	13.62	19.91	14.33	56.56	10.73	6.525	21.59
Cohort FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
Region FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
Obs.	202,196	120,775	81,421	193,221	113,431	52,294	251,674	143,733	107,941

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

Table 7, presents 2SLS estimates of equations 1 and 5 for the per capita consumption aggregate and the weighted per capita consumption aggregate constructed from the LSMS data. These two indexes give very close results. However, values of the F-test of overidentification restrictions are low except for two sub-groups (when I use instrument  $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$ ). It is around 10 for sample T3 when both rural and urban are taken together. The estimate range is wide and goes from 8.3 percent to 9.5 percent. For the urban T3 sample, the point estimate is from 12.6 to 13.2 percent.



Table 7: Estimations of the effect on education on LSMS consumption aggregates)

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)
	All	T1 Rural	Urban	All	T2 Rural	Urban	All	T3 Rural	Urban
LSMS- Deaton and Zaidi Consumption aggregate (2002)									
OLS	0.0695*** (0.00450)	0.0540*** (0.00545)	0.0774*** (0.00795)	0.0692*** (0.00527)	0.0556*** (0.00568)	0.0797*** (0.00762)	0.0727*** (0.00550)	0.0543*** (0.00502)	0.0802*** (0.00524)
R2	0.484	0.343	0.459	0.429	0.341	0.427	0.461	0.357	0.415
IV: $S_{j58} * T$	-0.0961 (0.183)	0.179 (0.654)	0.0143 (0.0538)	0.0374 (0.121)	0.529 (1.537)	0.135 (0.107)	0.0283 (0.0394)	-0.0811 (0.139)	0.110* (0.0578)
R2	-0.276	-0.223	0.209	0.332	-5.974	0.254	0.332	-0.167	0.468
F-test	1.225	0.0739	2.609	3.442	0.0877	2.674	12.12	2.040	2.562
IV: $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$	-0.113* (0.0650)	-0.0192 (0.0521)	0.0378 (0.0508)	0.0544 (0.0687)	0.163 (0.101)	0.138 (0.116)	0.0950*** (0.0253)	-0.00353 (0.0822)	0.132*** (0.0216)
R2	-0.412	0.080	0.294	0.349	-0.069	0.244	0.362	0.169	0.252
F-test	3.348	3.979	1.277	7.318	2.028	1.259	9.701	2.160	21.68
LSMS- Deaton and Zaidi Consumption weighted aggregate (2002)									
OLS	0.0642*** (0.00443)	0.0514*** (0.00510)	0.0730*** (0.00777)	0.0654*** (0.00516)	0.0527*** (0.00587)	0.0754*** (0.00837)	0.0688*** (0.00540)	0.0512*** (0.00528)	0.0768*** (0.00556)
R2	0.409	0.292	0.420	0.383	0.295	0.381	0.426	0.304	0.386
IV: $S_{j58} * T$	-0.124 (0.229)	0.0266 (0.473)	-0.0260 (0.0627)	0.0379 (0.101)	0.500 (1.457)	0.0866 (0.0702)	0.0379 (0.0388)	-0.0820 (0.136)	0.0916 (0.0559)
R2	-0.611	0.171	-0.056	0.287	-6.200	0.292	0.310	-0.336	0.298
F-test	1.225	0.0739	2.609	3.442	0.0877	2.674	12.12	2.040	2.562
$\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$	-0.0985 (0.0636)	-0.0209 (0.0537)	0.0132 (0.0466)	0.0535 (0.0583)	0.140 (0.0860)	0.0995 (0.0852)	0.0834*** (0.0270)	0.00686 (0.0629)	0.126*** (0.0228)
R2	-0.379	0.019	0.173	0.302	-0.047	0.276	0.328	0.140	0.223
F-test	3.353	4.165	1.295	7.209	2.028	1.259	9.694	2.125	21.47
Cohort FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
District FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
Obs.	2,914	1,499	666	2,779	1,403	636	3,481	1,639	902

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

The lower values of the F-test for the LSMS data may be explained by smaller sample size of sub-samples. Thus, if few observations are available by region and year of birth,  $C_{ijt}$  and  $S_{ijt}$  are not measured with accuracy. Thereafter, I focus on results from this IPUMS data.

## 5.4 Quality bias

Interpretation of these IV estimates may be delicate if the UPE program affects both the quantity and the quality of education Duffo [2000]. If that were to be true, the UPE program's effect on consumption would illustrate these two effects. To test whether the UPE program affects the quality of education, I compare returns from T1, T2 and T3 that have been differently affected by the UPE program. Individuals from T1 have not been exposed to the program since the beginning of their education. From T1 to T2, schools had to welcome an increasing number of pupils and quality of education may have been lowered by those rapid changes. T3 has been indirectly affected by the UPE program. It has also been exposed to structural changes that started in 1986 when Tanzania signed agreements with the IMF and the World Bank and when quality of education was defined as the new priority (Bonini, 2003). By focusing on the second line of table 6, one notes that IV estimates slightly fluctuate between T1, T2 and T3. In urban areas, these coefficients are not statically different. Despite the massive education program and structural changes, this result suggests that schools managed to maintain quality of education. However, coefficients vary more in rural areas: the points estimates decrease from T1 to T2 and increase from T2 to T3. This evolution may reflect changes in quality. When I use  $T * S_{j1945}$  as instrument, estimates of returns to education moves from 6.7 percent (T1) to 4.9 percent (T2). For the 2-SLS estimate with instrument  $\sum_{t=1961}^{1978} \gamma_t * S_{j1945}$ , the evolution is similar and the point estimate decreases by 1.9 percent. However, coefficients from T1 to T2 are not statically different.

It comes to the conclusion that the UPE program did not affect quality of education in urban areas. Regarding rural areas, the universalization of primary education generated a drop in quality but that it is not statically significant. Thus, it supports the idea that the UPE program affects consumption mainly through the quantity of education. The same interpretation could not be made between T2 and T3: estimate for T3 is statically different from T2. T3 coefficients are probably higher because they illustrate both the quantity impact of the UPE program and the quality impact generated by subsequent structural changes.

## 5.5 Returns to education by sector of activity

So far, returns to education have been estimated for the whole population and for rural and urban areas. However, they can vary from one sector of activity to another. In this section, I turn my attention to this purpose and I consider the endogeneity of education (the sample selection issue is breached later in section 7). First line of table 8 presents OLS results. These results are consistent with the literature: returns to education are lower in agriculture (4 percent) than in the self-employment sector (from

7.1 to 7.3 according to the subsample) and in the formal sector (from 6 to 6.4 percent). 2SLS estimations without sample selection correction are reported line 2. when I focus on treated groups T1 and T2, the estimate range is from 5.7 to 7.2 percent for the agricultural sector, from 7.5 to 8.8 percent for the self-employment sector, and from from 6.6 to 9.1 percent for the formal sector. Except for the self-employment sector, 2SLS estimates tend to give higher estimates than OLS estimates. Regarding T3 subsample, the returns to education are about 10 percent in each sector. These higher estimates may come from structural changes that improved quality of education.

Table 8: Returns to education by sector of activity

Sector	agri (1)	self (2)	formal (3)	agri (4)	self (5)	formal (6)	agri (7)	self (8)	formal (9)
OLS estimates									
	0.0400*** (0.00280)	0.0714*** (0.00426)	0.0639*** (0.00408)	0.0421*** (0.00265)	0.0744*** (0.00487)	0.0626*** (0.00364)	0.0422*** (0.00256)	0.0733*** (0.00398)	0.0602*** (0.00371)
R2	0.219	0.477	0.320	0.219	0.470	0.316	0.206	0.460	0.305
IV estimates									
	0.0720** (0.0328)	0.0882*** (0.0275)	0.0655*** (0.0221)	0.0574*** (0.0198)	0.0746** (0.0316)	0.0914*** (0.0230)	0.0979*** (0.0200)	0.107*** (0.0229)	0.115*** (0.0431)
R2	0.024	0.362	0.243	0.044	0.358	0.224	-0.015	0.345	0.169
F-test	12.30	16.74	7.615	21.88	12.46	19.58	10.34	26.86	8.615
IV estimates with sample selection correction (IV for occupation equation: $T * S_{j1945}$ )									
	0.0730** (0.0325)	0.0893*** (0.0280)	0.0655*** (0.0219)	0.0586*** (0.0194)	0.0760** (0.0323)	.0919*** (0.0228)	0.0999*** (0.0199)	0.110*** (0.0226)	0.116*** (0.0428)
mills2	-0.00735 (0.00473)			-0.00754** (0.00334)			-0.0135** (0.00524)		
mills3		0.00166 (0.00195)			0.000970 (0.00280)			0.00630* (0.00361)	
mills4			0.00312 (0.0246)			0.00848 (0.0195)			0.0135 (0.0175)
R2	0.023	0.362	0.243	0.044	0.358	0.011	-0.020	0.343	0.167
F-test	12.31	16.57	7.659	21.83	12.28	25.014	9.871	27.27	8.911
IV estimates with sample selection correction (IV for occupation equation: $T * S_{j1945}$ and $N_5$ )									
	0.0733** (0.0325)	0.0885*** (0.0282)	0.0645*** (0.0219)	0.0594*** (0.0193)	0.0755** (0.0324)	.091*** (0.0226)	0.103*** (0.0195)	0.109*** (0.0228)	0.115*** (0.0427)
mills2	-0.0103** (0.00459)			-0.0125*** (0.00374)			-0.0321*** (0.00599)		
mills3		0.000541 (0.00213)			0.000126 (0.00299)			0.00532 (0.00375)	
mills4			0.0203 (0.0219)			0.0219 (0.0173)			0.0349** (0.0141)
R2	0.023	0.362	0.244	0.043	0.358	-0.808	-0.027	0.343	0.171
F-test	12.33	16.49	7.438	21.80	12.26	24.72	9.749	26.92	8.635
Cohort FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
Region FE	YES	YES	YES	YES	YES	YES	YES	YES	YES
Obs.	136,773	31,298	32,276	128,686	32,951	29,864	162,278	50,097	36,806

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

## 6 The effect of education on mobility

Education increases potential welfare but can also be regarded as a mobility factor. While Tanzania is an agriculture-based country where 62 percent of the Tanzanian workforce is in agriculture, it could be interesting, both from a macroeconomic and a microeconomic perspective, to identify the effect of education on the distribution among areas and among sectors of activity.

### 6.1 Mobility across rural and urban areas

I define  $M_{ijt}$  the dummy variable taking value 1 for individuals that migrated to another region and went to urban areas and I estimate the following equation:

$$M_{ijt} = \alpha + \beta_j + \beta_t + \rho S_{ijt} + \delta T * X_{j,1945} + \epsilon_{ijt} \quad (6)$$

Table 9 gathers the results. OLS estimates from the IPUMS data show that one additional year of education increases the probability of migrating to urban areas in another region by 0.3 percent. Even if the impact is low, coefficients are significant for all samples. IV estimates lead to opposite results: coefficients range is from  $-0.8$  percent to  $-1.5$  percent. This change of sign is probably due to a specific effect of the treatment. By introducing agricultural classes in the curriculum, individuals who have been to school during the program may have been more likely to stay in the agricultural sector in rural areas.

However, this interpretation is limited because  $M_{ijt}$  does not take into account migration within regions.<sup>13</sup> Thus, I turn to the LSMS data and estimate equation (6) by first considering only migration between regions and second, by taking into account migration between districts. The two OLS estimates give similar results : the effect of education on the probability to move in urban areas in another region is close to 0 and is never significant. Thus, introducing migration between districts does not seem to change the results. However, while LSMS estimations capture more detailed migration flows, it still does not include migration within the same district. 2SLS LSMS estimates show an identical pattern to IPUMS estimates but they are difficult to interpret because the F-test range is from 4.2 to 6.3.

---

<sup>13</sup>IPUMS data do not precise the district of birth.

Table 9: Marginal average effect of education on the probability to migrate to urban areas

Sample	T1 (1)	T2 (2)	T3 (3)
IPUMS data when $M_{ijt}$ = migration between region			
OLS	0.00174*** (0.000240)	0.00184*** (0.000271)	0.00186*** (0.000334)
R2	0.61	0.61	0.59
IV	-0.0111*** (0.00215)	-0.00874*** (0.00236)	-0.0152*** (0.00308)
C.F.	.0128*** (.0023)	.0106*** (.0025)	.0171*** (.0032)
R2	0.61	0.59	0.59
F-test	36.02	60.98	39.35
Observations	202,196	193,221	251,674
Observations	202,196	193,221	251,674
LSMS data when $M_{ijt}$ = migration between region			
OLS	0.000302 (0.000719)	0.000220 (0.00315)	0.000720 (0.00307)
R-squared	0.510	0.486	0.495
IV	-0.00492 (0.0171)	-0.00273 (0.0474)	-0.00212 (0.00848)
C.F	.619** (0.2894)	.359 (0.2231)	.146 (0.2169)
R-squared	0.071	0.032	0.208
F-test	4.247	6.295	4.735
LSMS data when $M_{ijt}$ = migration between region and district			
OLS	0.000729 (0.00213)	0.000792 (0.00700)	0.00157 (0.00454)
R-squared	0.466	0.435	0.443
IV	-0.0186 (0.0540)	-0.00743 (0.0759)	-0.0235 (0.0683)
C.F	.717** (0.2266)	.325* (0.1727)	.538 (0.1888)
R-squared	-0.539	-0.258	-0.846
F-test	4.247	6.295	4.735
Observations	1,468	1,463	1,792
Cohort FE	YES	YES	YES
region FE	YES	YES	YES

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity and dummies for the destination region.

## 6.2 Effect of education on the labor market participation

I estimate the multinomial model where  $A_{ijtk}$  is the occupational choice taking value 1 if individuals do not work or are unpaid, 2 if individuals work in the agricultural sector, 3 if individuals work in the self-employment sector and 4 if they work in the formal sector. The probabilities have the following functional form :

$$P[A_{ij t} = k] = \frac{\exp(\beta_{jk} + \beta_{kt} + \gamma_k S_{ij t} + \delta_k T * X_j)}{1 + \exp(1 + \sum_{l=1}^4 (\beta_{jl} + \beta_{lt} + \gamma_l S_{ij t} + \delta_l T * X_j))} \quad (7)$$

Results are produced with OLS and with 2SLS estimates. To get rid of endogeneity issues, I instrument education and I follow a two-step Control Function approach [Wooldridge, 2014]. After obtaining the predicted residual from equation (1), I insert it in the multinomial model. The CF term added can be used to test whether education is endogenous. Table 10 reports the average marginal effects of education for T1, T2 and T3 samples and the coefficients of the CF term. OLS estimates show that one additional year of education decreases the probability of working in the agricultural sector from 0.3 percent to 0.4 percent and decreases the probability of working in the self-employment sector from 0.2 percent to 0.3 percent. On the opposite, it increases the probability of working in the formal sector by 0.7 percent. These effects are relatively small. Instrumenting education leads to different conclusions. It suggests that education has no effect on the probability of working in the formal sector. It decreases the probability of working in the self-employment sector and increases the probability of working in agriculture. However, these results may be specific to the treatment since the UPE program supported agriculture. One has to highlight that the CF terms are significant in the self-employment and the formal sectors for both T1 and T2. Thus, education is endogenous for these sub-samples: correlation between education and abilities seems to be higher for people working in the formal sector or the self-employment sector.

Table 10: Average marginal effect of education on the probability of working in each sector of activity

Sample	OLS				IV: $S_j, 1945 * T$			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Activity	Don't paid Don't work	agri.	self employed	formal	Don't paid Don't work	agri	self employed	formal
T1	-0.000654*** (9.34e-05)	-0.00466*** (0.00123)	-0.00178* (0.00106)	0.00709*** (0.000349)	0.00183 (0.00178)	0.0113** (0.00545)	-0.0124*** (0.00413)	-0.000748 (0.00324)
C.F.						0.183 (0.1852)	0.436* (0.1749)	0.485** (0.1802)
F-test					36.02	36.02	36.02	36.02
Obs.	203496	203496	203496	203496	203496	203496	203496	203496
T2	-0.000542*** (0.000109)	-0.00408*** (0.00151)	-0.00273** (0.00139)	0.00735*** (0.000403)	-0.000529 (0.00123)	0.00998*** (0.00282)	-0.0107*** (0.00317)	0.00125 (0.00197)
C.F.						-0.0593 (0.1269)	0.141 (0.1377)	0.180 (0.1182)
F-test					60.98	60.98	60.98	60.98
Obs.	194419	194419	194419	194419	194419	194419	194419	194419
T3	-0.000783*** (8.21e-05)	-0.00298** (0.00137)	-0.00329** (0.00134)	0.00706*** (0.000310)	0.000564 (0.000755)	0.0114*** (0.00397)	-0.0136** (0.00587)	0.00169 (0.00315)
C.F.						0.0585 (0.0778)	0.278** (0.0860)	0.285** (0.0942)
F-test					39.45	39.45	39.45	39.45
Obs.	253288	253288	253288	253288	253288	253288	253288	253288
F-test						39.45		
Obs.		253288				253288		
Cohort FE	YES	YES	YES	YES	YES	YES	YES	YES
Region FE	YES	YES	YES	YES	YES	YES	YES	YES

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. C.F. terms reported are coefficients from the multinomial model. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.



## 7 Sample selection correction

Estimates of the returns to education are biased in case of sample selection. Duflo [2000] adopts two methods to deal with this selection issue and to measure the returns to education for individuals working for a wage. Her first strategy is to control whether the returns to education are stable when she takes into consideration the entire population and imputes an income for self-employed individuals. In this study, I should not encounter this sample selection issue since consumption is available for all households. However, this problem has to be addressed when returns to education are estimated for non random sub-samples. It happens when the choice of being in a sub-group is affected by education.

Results of table 6 are reassuring: it shows that returns to education for the whole sample are between returns to education in rural areas and returns to education in urban areas. Differences between the whole sample and rural or urban areas are never significant. Similarly, the average returns to education is between the returns to education by sector (see table 8).

The second way of identifying whether results are biased by sample selection is to test if education influences the probability of moving in a specific area and in that case, add a correction term to equation (3).

### 7.1 Rural and Urban returns to education with sample selection correction

To examine whether the returns to education in rural and urban areas are biased by sample selection, I refer to the model discussed by Wooldridge [2010] that considers sample selection model when one of the explanatory variable is endogenous. The first stage consists of obtaining the Inverse Mills Ratio (IMR) from the 2SLS estimate of the selection equation (equation 6). Then, I add the IMR in the principal equation (equation 3) and I estimate it by 2SLS. Results are presented in Table 11 and reflect that the sample selection bias is negligible. Comparison of the first and the second line shows that results are not sensitive to the introduction of the IMR. IMR coefficients are close to 0 but are always significant (except in rural areas for sample T3).

Table 11: IV estimates of returns to education by areas

Instrument	rural (1)	urban (2)	rural (3)	urban (4)	rural (5)	urban (6)
IV estimate with instrument $\gamma_t * S_{j1945}$						
	0.0614*	0.0942***	0.0415*	0.0851***	0.102***	0.0969***
	(0.0373)	(0.0140)	(0.0243)	(0.0202)	(0.0258)	(0.0139)
R2	0.091	0.254	0.105	0.197	0.006	0.218
F-test	10.74	13.62	14.33	56.56	6.525	21.59
IV estimate with instrument $\gamma_t * S_{j1945}$ and Heckman correction						
	0.0638*	0.0922***	0.0443*	0.0872***	0.104***	0.0955***
	(0.0344)	(0.0139)	(0.0233)	(0.0143)	(0.0257)	(0.0145)
IMR	-0.000243***	0.000330***	-0.000182**	0.000291**	-8.83e-05	0.000252**
	(7.67e-05)	(8.23e-05)	(8.41e-05)	(0.000122)	(0.000103)	(9.96e-05)
R-squared	0.088	0.260	0.105	0.261	-0.000	0.222
F-test	10.69	45.70	14.55	41.32	6.198	21.65
Cohort FE	YES	YES	YES	YES	YES	YES
Region FE	YES	YES	YES	YES	YES	YES
Observations	120,775	81,421	113,431	79,790	143,733	107,941

Note: Source: the IPUMS data. Standard errors are clustered at the birth region level and are reported in parentheses. \*\*\*, \*\*, \* means respectively that the coefficient is significantly different from 0 at the level of 1%, 5% and 10%. C.F. terms reported are coefficients from the multinomial model. Additional controls are the population aged 7 to 13 in 1958, the percentage of people living in rural areas in 1958, the household's size and the sector of activity.

## 7.2 Returns to education by sector of activity with sample selection correction

To correct the sample selection that may bias the results of the returns to education by sector of activity, I also adopt the two-stage model proposed by Wooldridge [2010]. For the first stage, I regress the occupational choice variable  $A_{ijt}$  on the instrument  $S_{jt} * T$ . By using the predicted probabilities from this model, I compute the inverse Mill's ratios for each outcome. Then, I include them in the respective consumption equations which are estimated with 2SLS. As Wooldridge [2010] underlines, if the same instruments are used for the occupational equation and for the consumption equation, the introduction of the Mills ratio generates collinearity that can affect performance in the case small samples. For not suffering from this issue, it is preferable to include another instrument in the occupational equation that is not present in the consumption equation. Even if the sub-samples' size are very large, I test whether adding another instrument (the number of children younger than 5  $N_5$ ) in the occupational equation changes the results. Results with sample selection correction are reported at the end of table 8. There is little evidence of sample selection bias: coefficients of the Mill's ratio are close to 0 and are not statically significant except for the agricultural sector. Further, the introduction of sample selection corrections do not statically change the 2SLS estimates.

This interpretation is still valid when I add another instrument in the

occupational equation: coefficients do not change and F-test is not improved.

## 8 Conclusion

This paper brings to light that the massive primary education program implemented in Tanzania contributed to reduce inequalities among regions. It tended to equalize the access to schools by developing education in the most deprived areas. The effect of this program did not suddenly stop with the end of the program, but on the contrary, persisted for the next age-cohorts. In this respect, the Tanzanian government fulfilled its goal by equalizing the access of basic education.

This increase in education also resulted into changes in the distribution of individuals among sectors of activities and among locations; the government managed to encourage people to work in the agricultural sector instead of the self-employment sector. This effect is probably related to the specific nature of the UPE program since a significant part of the school curriculum were devoted to agricultural classes.

When I instrument education by using variation in intensity of the UPE program, I find that education increases household's consumption aggregates between 7.4 percent and 10.5 percent. In urban areas, people who have been partially treated have similar returns that fully treated people. However, this is no longer the case in rural areas : returns drop for fully treated people probably because quality of schools is affected by the massive UPE program. It is not surprising since classrooms were likely to be overcrowded with the lack of schools and the abrupt increase of enrolled pupils. However, returns to education increase again for the next age-cohorts suggesting that schools managed to adapt the schooling system.

This paper underlines that the effect of education on household's consumption is heterogenous according to the sector of activity of the household's head and to the household's location. The education effect is larger in urban areas, in the self-employment sector and in the formal sector than in rural areas and in the agricultural sector. Thus, this paper entails to give an overview of the distribution of the returns to education between sectors and contributes to the scarce literature on the returns to education in agriculture in developing countries.

Instead of looking at the effect of education on the probability of working in each sector, it would be interesting to study the effect of education on the probability of cumulating several jobs. Indeed, if education gives access to a larger range of jobs, it could constitute a coping strategy in order to protect against productivity shocks that make vulnerable people working in the agricultural sector.

## References

- Nathalie Bonini. Un siècle d'éducation scolaire en tanzanie. *Cahiers d'études africaines*, (1):40–62, 2003.
- David Card. Estimating the return to schooling: Progress on some persistent econometric problems. *Econometrica*, 69(5):1127–1160, 2001.
- Angus Deaton and Salman Zaidi. *Guidelines for constructing consumption aggregates for welfare analysis*, volume 135. World Bank Publications, 2002.
- Esther Duflo. Schooling and labor market consequences of school construction in indonesia: Evidence from an unusual policy experiment. Technical report, National Bureau of Economic Research, 2000.
- Zvi Griliches. Research expenditures, education, and the aggregate agricultural production function. *The American Economic Review*, 54(6): 961–974, 1964.
- Kabiru Kinyanjui et al. Development policy and educational opportunity: the experience of kenya and tanzania. 1980.
- Marlaine E Lockheed, T Jamison, and Lawrence J Lau. Farmer education and farm efficiency: A survey. *Economic development and cultural change*, 29(1):37–76, 1980.
- John Maluccio. Endogeneity of schooling in the wage function: Evidence from the rural philippines. *Food Consumption and Nutrition Division Discussion Paper*, 54, 1998.
- Paul S Maro and Wilfred FI Mlay. Decentralization and the organization of space in tanzania. *Africa*, 49(03):291–301, 1979.
- Denis Martin. *Tanzanie: l'invention d'une culture politique*. KARTHALA Editions, 1988.
- DJ McKenzie. Measure inequality with asset indicators. cambridge, ma: Bureau for research and economic analysis of development. *Center for International Development, Harvard University*, 2003.
- Julius K Nyerere. Ujamaa: the basis of african socialism. *The Journal of Pan-African Studies*, 1:4–11, 1987.
- Ricardo Sabates, Jo Westbrook, and Jimena Hernandez-Fernandez. *The Health and Education Benefits of Universal Primary Education for the Next Generation: Evidence from Tanzania. CREATE Pathways to Access. Research Monograph No. 62*. ERIC, 2011.

David E. Sahn and David Stifel. Exploring alternative measures of welfare in the absence of expenditure data. *Review of Income and Wealth*, 49(4): 463–489, 2003. ISSN 1475-4991. doi: 10.1111/j.0034-6586.2003.00100.x. URL <http://dx.doi.org/10.1111/j.0034-6586.2003.00100.x>.

Seema Vyas and Lilani Kumaranayake. Constructing socio-economic status indices: how to use principal components analysis. *Health policy and planning*, 21(6):459–468, 2006.

Jeffrey M Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.

Jeffrey M Wooldridge. Quasi-maximum likelihood estimation and testing for nonlinear models with endogenous explanatory variables. *Journal of Econometrics*, 182(1):226–234, 2014.

# Appendices

Table 12: Descriptive statistics

Variable	Obs.	Mean	Std. Dev.	Min	Max
IPUMS data					
Ln (index_y)	931066	1,044	0,862	0	2,644
Education	935456	5,454	3,535	0	15
Primary completion	938026	0,639	0,480	0	1
Edu. level in 1958	938026	3,006	0,829	1,631	5,829
Children younger than 5	938026	0,657	0,843	0	9
HH size	938026	4,832	3,068	1	30
Urban	938026	1,429	0,495	1	2
Men	938026	0,494	0,500	0	1
formal sector employer	789366	0,127	0,333	0	1
self employment	789366	0,179	0,383	0	1
unpaid agri	789366	0,015	0,120	0	1
percentage rural in 1958	938026	0,682	0,087	0,161	0,959
ln (school size in 19858)	938026	8,706	0,492	6,234	9,181
LSMS data					
Consumption	6823	783 477,900	2029741	16 792,350	104 000 000,000
Weighted consumption	6823	1078339	3328438	28 366,140	180 000 000,000
Education	6740	6,349	4,170	0	20
Age of the HH head	6823	45,367	9,953	30	67
Gender of the HH head	6823	0,235	0,424	0	1
Size of the HH	6823	5,591	3,052	1	55

Sources: The 2002 census (IPUMS data) and the three pooled waves of the LSMS data.

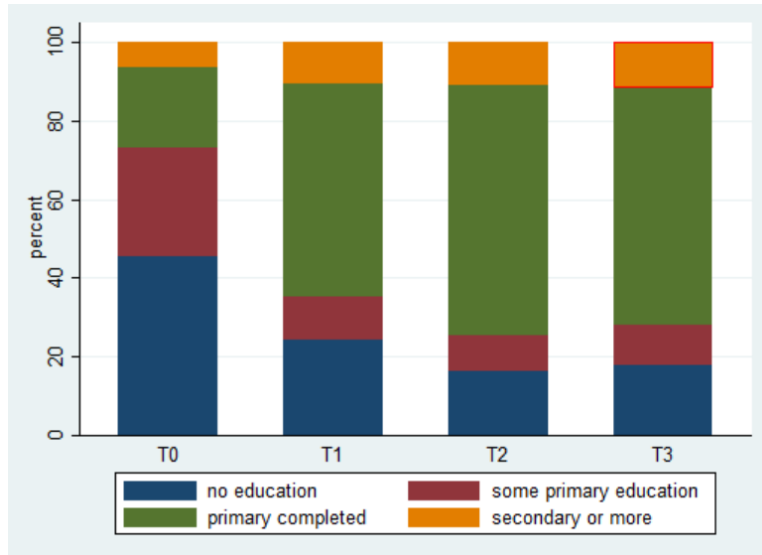


Figure 3: Evolution of the education attainment by Age-Cohort.

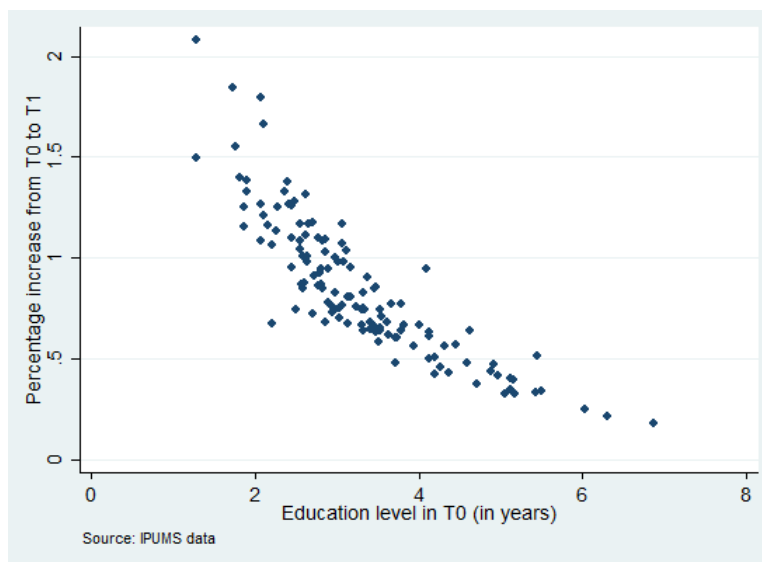


Figure 4: Evolution of the education attainment by district according to the education level in 1967.

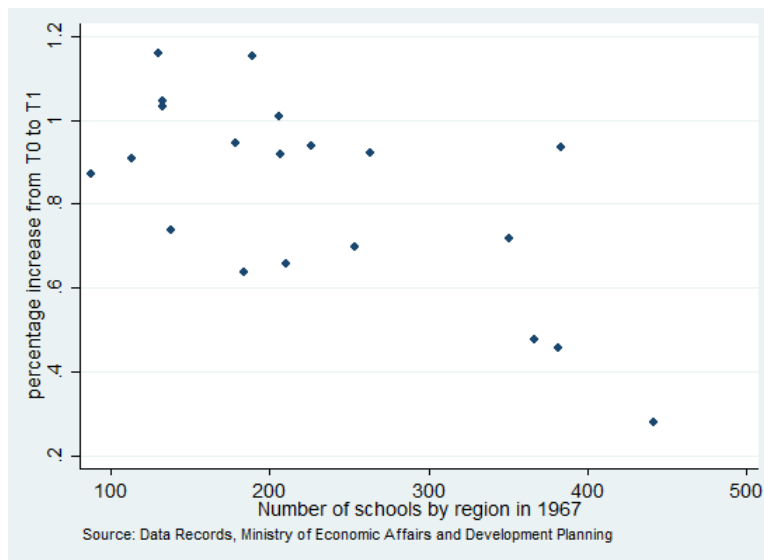


Figure 5: Evolution of the education attainment by region according to the number of schools in 1967.