



**Working papers series**

---

**WP ECON 18.03**

***Estimation of competing risks duration  
models with unobserved heterogeneity  
using hsmlogit***

David Troncoso-Ponce

Universidad Pablo de Olavide

**Keywords:** Duration analysis, Unobserved heterogeneity, d2 ml method, hshaz, hshaz2, hsmlogit, Hessian matrix, Multinomial Logit, Competing risks, Stata.

**JEL Classification:** C23, C25, C41, C54, C55, J64, J68.



**Department of Economics**

---

# Estimation of competing risks duration models with unobserved heterogeneity using `hsmlogit`

David Troncoso Ponce  
Pablo de Olavide University  
Seville, Spain  
Email: dtropon@upo.es

## Abstract.

This article presents `hsmlogit`, a new Stata command that estimates multispells discrete time competing risks duration models with unobserved heterogeneity. `hsmlogit` allows for the estimation of one, two and up to three competing risks, as well as a maximum of five points of support for the identification of unobserved heterogeneity distribution ([Heckman and Singer, 1984]). The main contribution of `hsmlogit` is that allows for exploiting the richness of large longitudinal micro datasets, by estimating competing risks duration models, instead of one-risk models (such as `hshaz` and `hshaz2`), as well as it takes into account the presence of unobserved heterogeneity affecting transition rates. In addition to this, and taking into account the larger size of longitudinal micro datasets used for the estimation of discrete time duration models, `hsmlogit` also provides the algebraic expressions of both first and second order derivatives that, respectively, define the gradient vector and Hessian matrix, which significantly reduce time required to achieve model convergence.

**Keywords:** Duration analysis, Unobserved heterogeneity, `d2 ml` method, `hshaz`, `hshaz2`, `hsmlogit`, Hessian matrix, Multinomial Logit, Competing risks models

## Acknowledgement

I am grateful to professor Stephen Jenkins for his helpful comments and suggestions that have contributed to significantly improve this work, and for allowing me to use Stata code of his `hshaz`'s command syntax. I also thank financial support of research project SEJ-6882 from Junta de Andalucía.

## 1 Introduction

Empirical studies on individual decisions have experienced an important increase in recent years due to the boost of large and rich longitudinal micro datasets put available to the research community. Specially for the field of empirical Labor Economics focused on the estimation of labor market transition rates, the recent availability of large longitudinal micro datasets allows for capturing the presence of unobserved heterogeneity (UH, hereafter) components that affect the estimated transition rates. However, an important number of

these empirical studies that incorporates the presence of UH mainly focused on one-risk duration models,<sup>1</sup> that analyze transition rates towards an only destination (for example, transitions from employment to unemployment, ignoring the existence of other destinations, such as inactivity, or finding another job). This article presents `hsmlogit`, a new Stata command that estimates multi-spells discrete time competing risks duration models with UH. `hsmlogit` allows for the estimation of one, two and up to three competing risks, as well as a maximum of five points of support for the identification of unobserved heterogeneity distribution ([Heckman and Singer, 1984]).

The main contribution of `hsmlogit` is that allows for exploiting the richness of large longitudinal micro datasets, by estimating competing risks duration models, instead of one-risk models (such as `hshaz` and `hshaz2`), as well as it takes into account the presence of unobserved heterogeneity affecting transition rates. In addition to this, and taking into account the larger size of longitudinal micro datasets used for the estimation of discrete time duration models, `hsmlogit` also provides the algebraic expressions of both first and second order derivatives that define the gradient vector and Hessian matrix, respectively, which significantly reduce time required to achieve model convergence [Gould et al., 2010].

The rest of the article structures as follows: Section 2 describes the longitudinal database used to obtain estimation results; the econometric model and `hsmlogit` command syntax are explained, respectively, in Sections 3 and 4; Section 5 presents estimation results, and Section 6 shows the advantages of providing the algebraic expressions of both the gradient vector and Hessian matrix. Finally, Section 7 concludes.

## 2 Database: The Continuous Sample of Working Histories

I analyze a longitudinal sample of workers in the Spanish labor market that comes from the *Continuous Sample of Working Histories* database (CSWH, hereafter). The CSWH is a longitudinal database that provides the working histories records of more than one million people, who represent a 4% non-stratified random draw from a target population, composed of any person with a contribution relation with the Spanish Social Security Administration. It includes both wage workers and recipients of Social Security benefits, namely, unemployment benefits, disability, survivor pension and maternity leave.<sup>2</sup>

<sup>1</sup>An exception is the work presented in [Troncoso-Ponce, 2016], where a two-states multi-spells discrete time competing risks duration model with UH is estimated to analyze the effect of apprenticeship contracts in the Spanish labor market.

<sup>2</sup>[García-Pérez, 2008], [Lapuerta, 2010], [Arranz and García-Serrano, 2011] and [Arranz, García-Serrano and Hernanz, 2013] contain a deep exposition about features of CSWH as well as all necessary techniques to perform a duration analysis using working lives information.

The CSWH contains detailed information on each employment and unemployment episodes experienced by workers through their entire working histories. The information provided by the CSWH can be grouped into several categories: First, personal characteristics of workers (gender, age, nationality, educational level, residence place, and other personal characteristics). Second, job characteristics (type of labor contract, part-time coefficient, qualification level, and other job characteristics). Third, information on the employer (firm size, activity sector, and other firm characteristics). Furthermore, an important feature of the CSWH is that provides the beginning and termination dates of all employment and unemployment episodes, which takes special interest for duration analysis.

The estimation sample is composed of 48,246 low-educated and low-qualified young workers in the Spanish labor market for the period 2000-2014. The average age is 22.5 years-old, and 75.5% of them are males. The average number of employment episodes per worker is 8.9, lasting, on average, 7.13 months. Indeed, more than 25% of all employment episodes last 2 months or less, and only 5% last at least 24 months, which highlights the high turnover rate experienced by these workers. Multispell estimation sample has 1,316,611 observations. A **describe** output is shown below, describing the estimation sample's full varlist, as well as a **summarize** output to show the main descriptive statistics of the estimation sample's varlist.

## Description of the estimation sample's varlist

```
. describe , fullnames
      obs:      1,316,611
      vars:           59
      size:    127,711,267
```

29 Sep 2017 17:04

variable name	storage type	display format	value label	variable label
codind	long	%12.0g		id of individual
spell	int	%9.0g		Sequential number of the employment episode
j	int	%9.0g		Month of employment spell
exit3	byte	%8.0g	exit3	Exit from employment state (3 competing risks)
exit1	byte	%8.0g	exit1	Exit from employment state (1 only risk)
exit2	byte	%8.0g	exit2	Exit from employment state (2 competing risks)
cf	byte	%8.0g		Apprenticeship contract (=1)
ct	byte	%8.0g		Temporary contract (=1)
lnjemp	float	%9.0g		Log(t)
lnjemp2	float	%9.0g		Log(t)^2
lnjemp3	float	%9.0g		Log(t)^3
month1	byte	%8.0g		Month 1
month2	byte	%8.0g		Month 2
month3	byte	%8.0g		Month 3
month6	byte	%8.0g		Month 6
month12	byte	%8.0g		Month 12
month18	byte	%8.0g		Month 18
month24	byte	%8.0g		Month 24
month36	byte	%8.0g		Month 36
month48	byte	%8.0g		Month 48
female	byte	%8.0g		Female (=1)
age16tv	byte	%9.0g		Current age - 16
age16tv2	int	%9.0g		(Current age - 16)^2
educcompul1	byte	%8.0g		Education: Compulsory stage #1
educcompul2	byte	%8.0g		Education: Compulsory stage #2
educless1	byte	%8.0g		Education: Less than compulsory stage #1
educless2	byte	%8.0g		Education: Less than compulsory stage #2
inmigra	byte	%8.0g		Not Spanish nationality (=1)
manufactory	byte	%8.0g		Economic sector: Manufacturing industry
highserv	byte	%8.0g		Economic sector: High qualified services
lowserv	byte	%8.0g		Economic sector: Low qualified services
comerce	byte	%8.0g		Economic sector: Commerce
highqualif	byte	%8.0g		Previous job: High qualification
midhighqualif	byte	%8.0g		Previous job: Mid-High qualification
midlowqualif	byte	%8.0g		Previous job: Mid-Low qualification
lowqualif	byte	%8.0g		Previous job: Low qualification
prevunemp	byte	%9.0g		Number of previous unemployment spells
prevtc	int	%9.0g		Number of previous temporary contracts
unrate	double	%10.0g		Quarterly regional unemployment rate (Q.r.u.r.)
unratexlnjemp	float	%9.0g		(Q.r.u.r.) x Log(t)
unratexlnjemp2	float	%9.0g		(Q.r.u.r.) x Log(t)^2
gremloyment	float	%9.0g		Quarterly employment growth rate (Q.e.g.r.)
gremloymentxlnjemp	float	%9.0g		(Q.e.g.r.) x Log(t)
gremloymentxlnjemp2	float	%9.0g		(Q.e.g.r.) x Log(t)^2
andal	byte	%8.0g		Spanish region: Andalucia
aragon	byte	%8.0g		Spanish region: Aragon
astur	byte	%8.0g		Spanish region: Asturias
balear	byte	%8.0g		Spanish region: Baleares

canar	byte	%8.0g	Spanish region: Canarias
cantab	byte	%8.0g	Spanish region: Cantabria
castman	byte	%8.0g	Spanish region: Castilla La Mancha
castleon	byte	%8.0g	Spanish region: Castilla Leon
valenc	byte	%8.0g	Spanish region: Valencia
extrem	byte	%8.0g	Spanish region: Extremadura
galic	byte	%8.0g	Spanish region: Galicia
murcia	byte	%8.0g	Spanish region: Murcia
navarr	byte	%8.0g	Spanish region: Navarra
vasco	byte	%8.0g	Spanish region: Pais Vasco
rioja	byte	%8.0g	Spanish region: La Rioja

---

Sorted by: codind spell j

## Descriptive statistics of the estimation sample's varlist

. sum codind spell j exit\* `varsaleE`

Variable	Obs	Mean	Std. Dev.	Min	Max
codind	1,316,611	3800765	2579758	1868	1.00e+07
spell	1,316,611	8.951002	13.46625	1	435
j	1,316,611	7.132385	9.23605	1	108
exit3	1,316,611	.3740976	.7186157	0	3
exit1	1,316,611	.2373564	.4254626	0	1
exit2	1,316,611	.369709	.7054998	0	2
cf	1,316,611	.16471	.3709188	0	1
ct	1,316,611	.83529	.3709188	0	1
lnjemp	1,316,611	1.407723	1.039737	0	4.682131
lnjemp2	1,316,611	3.062737	3.439172	0	21.92235
lnjemp3	1,316,611	7.66514	11.91313	0	102.6433
month1	1,316,611	.2119381	.4086814	0	1
month2	1,316,611	.1404713	.3474754	0	1
month3	1,316,611	.1057419	.3075072	0	1
month6	1,316,611	.054834	.227656	0	1
month12	1,316,611	.0209416	.1431891	0	1
month18	1,316,611	.0090034	.0944583	0	1
month24	1,316,611	.0052263	.072104	0	1
month36	1,316,611	.0014818	.0384661	0	1
month48	1,316,611	.0005301	.0230189	0	1
female	1,316,611	.2444549	.4297637	0	1
age16tv	1,316,611	5.573653	3.642791	0	22
age16tv2	1,316,611	44.33552	53.2066	0	484
educcompul1	1,316,611	.1649796	.3711623	0	1
educcompul2	1,316,611	.4394009	.4963143	0	1
educless1	1,316,611	.2096162	.4070349	0	1
educless2	1,316,611	.1860033	.3891095	0	1
inmigra	1,316,611	.1212887	.3264626	0	1
manufactory	1,316,611	.1599212	.3665331	0	1
highserv	1,316,611	.0580111	.2337645	0	1
lowserv	1,316,611	.1449327	.3520331	0	1
comerce	1,316,611	.1784703	.3829083	0	1
highqualif	1,316,611	.0061841	.0783953	0	1
midhighqua-f	1,316,611	.0349579	.1836734	0	1
midlowqualif	1,316,611	.3149009	.4644766	0	1
lowqualif	1,316,611	.6439571	.4788283	0	1
prevunemp	1,316,611	1.988224	2.428205	0	25
prevtc	1,316,611	5.621923	12.12324	0	431
unrate	1,316,611	11.85796	5.19547	3.9	36.87
unratexlnj-p	1,316,611	17.4539	17.60828	0	169.4221
unratexlnj-2	1,316,611	39.34931	58.8646	0	778.5148
greemployment	1,316,611	2.102568	3.801228	-14.00105	10.99764
greemploye-p	1,316,611	2.442792	7.083604	-62.03608	38.71074
greemploye-2	1,316,611	4.294824	19.12216	-274.8705	142.1837
andal	1,316,611	.2453321	.4302841	0	1
aragon	1,316,611	.0239395	.1528608	0	1

astur	1,316,611	.0198137	.1393599	0	1
balear	1,316,611	.0305026	.1719656	0	1
canar	1,316,611	.046025	.2095393	0	1
cantab	1,316,611	.013935	.1172214	0	1
castman	1,316,611	.0572325	.2322865	0	1
castleon	1,316,611	.0449579	.2072117	0	1
valenc	1,316,611	.1079005	.3102548	0	1
extrem	1,316,611	.0249443	.1559556	0	1
galic	1,316,611	.0732874	.2606078	0	1
murcia	1,316,611	.0322449	.1766499	0	1
navarr	1,316,611	.0085128	.0918711	0	1
vasco	1,316,611	.0262386	.1598441	0	1
rioja	1,316,611	.0053638	.0730411	0	1

### 3 Econometric model

This Section briefly describes the main features of the econometric models that will be estimated in Section 5. The main goal of this kind of models is to analyze duration spent by a population in a specific state (in this example, employment state), as well as to analyze the set of factors, observable and specially unobservable, that affect time spent in that state (see [Lancaster, 1992], [Allison, 1982] and [Jenkins, 1995]).

Let's consider an individual beginning an employment episode at time  $T = 1$  (time  $T$  is measured in month intervals). The worker is observed monthly during the employment episode until either he/she exits to another modeled state (such as, unemployment, or finding a new job), or the observation window ends (right censored observations). Employment duration is analyzed by estimating the hazard rate out of employment at each observed month. Depending on the number of exits (i.e. risks) modeled by the command's user, `hsmlogit` can estimate two different functional forms for the hazard rate.

Single-risk models use a *Logit* functional form to characterize the hazard rate, given by the following expression:

$$h(t|x, \eta) = \frac{\exp(\lambda(t) + x\beta + \eta)}{1 + \exp(\lambda(t) + x\beta + \eta)} \quad (1)$$

And competing risks models use a *Multinomial Logit* functional form to characterize the hazard rates:

$$h_d(t|x_d, \eta) = \frac{\exp(\lambda_d(t) + x_d\beta_d + \eta)}{1 + \sum_{d=1}^D \exp(\lambda_d(t) + x_d\beta_d + \eta)} \quad (2)$$

Assuming that  $h = \sum_{d=1}^D h_d$ , where  $d = 1, \dots, D$  and  $D = \{1, 2, 3\}$  depending on the total number of risks modeled by the command's user.



As the two expressions above show, the hazard rate at month  $T = t$  depends on time (months) spent in the current unemployment state (i.e. duration dependence), captured by  $\lambda(t)$ , as well as on a set of covariates summarized by  $x$  vector, that may contain both time-fixed and time-varying covariates. Furthermore, the hazard rate also depends on an unobserved component given by  $\eta$ , that measures factors, such as job search effort, job networking, motivation, ability, etc, that are unobserved to the researcher and may affect the transition rate out of employment.

For the case of one-risk models, the contribution to the likelihood function of an individual  $i$  is given by the following expression:

$$L_i = \sum_{j=1}^P \pi_j \left\{ \prod_{t=1}^{T_i} \frac{h(T=t|\lambda(t), x_{it}, \eta_j)^{y_{it}}}{(1 - h(T=t|\lambda(t), x_{it}, \eta_j))^{(1-y_{it})}} S(T=t|\lambda(t), x_{it}, \eta_j)^{(1-y_{it})} \right\} \quad (3)$$

Where dependent variable  $y_{it} = \{0, 1\}$  denotes a dummy variable that takes value 1 if worker  $i$  exits out from employment at month  $T = t$ , and takes value zero otherwise.<sup>3</sup> Expression given by  $h(T=t|\lambda(t), x_{it}, \eta_j)$  denotes the hazard rate observed at month  $T = t$ , and  $S(T=t|\lambda(t), x_{it}, \eta_j)$  denotes the survival rate observed at month  $T = t$ , that estimates the cumulative probability of being employed (from the month  $T = 1$ ) until the month  $T = t$ , and that is given by the following expression:

$$S(T=t|\lambda(t), x_{it}, \eta_j) = \prod_{s=1}^t (1 - h(T=s|\lambda(s), x_{is}, \eta_j)) \quad (4)$$

As expressions 1 and 4 show, the hazard rate observed at month  $T = t$  is conditional on the duration dependence  $\lambda(t)$  and on the set of covariates  $x_{it}$ . And the survival rate at month  $T = t$  is conditional on  $\lambda(s)$  and on the set of covariates  $x_{is}$ , observed at months  $s = 1, 2, \dots, t$ . Both the hazard and the survival rates also depend on belonging to the type of employed workers with unobserved characteristics given by  $\eta_j$ .<sup>4</sup>

The total likelihood function of single-risk models is given by:

$$L = \prod_{i=1}^N \sum_{j=1}^P \pi_j \left\{ \prod_{t=1}^{T_i} \frac{h(T=t|\lambda(t), x_{it}, \eta_j)^{y_{it}}}{(1 - h(T=t|\lambda(t), x_{it}, \eta_j))^{(1-y_{it})}} S(T=t|\lambda(t), x_{it}, \eta_j)^{(1-y_{it})} \right\} \quad (5)$$

`hsmlogit` command maximizes, using `d2 ml` method, the natural logarithm of  $L$  to estimate the model parameters.

<sup>3</sup>Dependent variable  $y_{it}$  refers to `dead(deadvar)` of `hsmlogit` command.

<sup>4</sup>It is assumed that unobserved characteristics do not vary with time and are not correlated to the rest of explanatory variables included in the specification of the hazard rate.

For the case of competing risks models, the contribution to the likelihood function of an individual  $i$  is given by the following expression:

$$L_i = \sum_{j=1}^P \pi_j \left\{ \prod_{t=1}^{T_i} \prod_{d=1}^D \{h_d(T_d = t | \lambda_d(t), x_{it}^d, \eta_j)^{y_{it}^d}\} S(T = t-1 | \lambda(t), x_{it}, \eta_j)^{(1 - \sum_{d=1}^D y_{it}^d)} \right\} \quad (6)$$

Where  $h_d(T_d = t | \lambda_d(t), x_{it}^d, \eta_j)$  denotes the hazard rate for the specific risk  $d = 1, \dots, D$  observed at month  $T = t$ , conditional on the duration dependence  $\lambda_d(t)$ , on the set of covariates  $x_{it}^d$ , and on belonging to the type of employed workers with unobserved characteristics given by  $\eta_j$ . Dependent variable  $y_{it}^d = \{0, 1\}$ , for  $d = 1, \dots, D$ , denotes a dummy variable that takes value 1 if worker  $i$  exits out from employment towards the destination  $d$  at month  $T_d = t$ , and takes value zero otherwise.

Unlike single-risk models, the survival function for competing risks takes into account the all possible risks faced by the individual at month  $T = t$ , and therefore takes the following expression:

$$S(T = t-1 | \lambda(t-1), x_{it-1}, \eta_j) = \prod_{s=1}^{t-1} \left( 1 - \sum_{d=1}^D h_d(T_d = s | \lambda_d(s), x_{is}^d, \eta_j) \right) \quad (7)$$

Similarly to single-risk models, the total likelihood function for competing risks is given by:

$$L = \prod_{i=1}^N \sum_{j=1}^P \pi_j \left\{ \prod_{t=1}^{T_i} \prod_{d=1}^D \{h_d(T_d = t | \lambda_d(t), x_{it}^d, \eta_j)^{y_{it}^d}\} S(T = t-1 | \lambda(t), x_{it}, \eta_j)^{(1 - \sum_{d=1}^D y_{it}^d)} \right\} \quad (8)$$

And likewise single-risk models, `hsmlogit` command maximizes, also using `d2 ml` method, the natural logarithm of  $L$  to estimate the model parameters for competing risks models.

### 3.1 The non-parametric identification of the UH distribution

Regarding the estimation of UH distribution, we assume the existence of unobserved factors affecting hazard rates, that if are ignored, may lead to spurious duration dependence, captured by  $\lambda(t)$  ([Van Den Berg, 2001]). A well known method to capture the effect of UH on the hazard rates is the proposed by [Heckman and Singer, 1984], by which the UH components are captured without imposing any parametric distribution function for the identification of UH distribution, but as a discrete mixture of several types of individuals with different values of UH components. Thus, it is assumed the presence of different types of workers who characterize themselves by having different levels of unobserved

variables (such as, ability, cognitive and non cognitive skills, social and networking capabilities, etc.), captured by the set of parameters  $\eta = \{\eta_1, \eta_2, \dots, \eta_P\}$ , that are estimated as regression's constant terms.<sup>5</sup> For each Type of worker  $j$ , characterized by  $\eta_j$ , an associated probability of being observed in the data, given by  $\pi = \{\pi_1, \pi_2, \dots, \pi_P\}$ , is also estimated jointly with the rest of the model parameters. Finally, the non-parametric discrete UH distribution is the result of the combination of these Types of workers, whose different values of UH are given by the vector  $\eta = \{\eta_1, \eta_2, \dots, \eta_P\}$  and by their associated probabilities  $\pi = \{\pi_1, \pi_2, \dots, \pi_P\}$ , are estimated jointly with the rest of the model parameters.

Furthermore, likewise **hshaz2** command, when more than two mass-points are specified by the command's user, **hsmlogit** also properly estimates mass-points probabilities using a *Multinomial Logit* function, rather than a *Logit* one, to compute the values of  $(\pi_1, \pi_2, \dots, \pi_P)$  (see [Troncoso-Ponce, 2017]). For example, when the UH distribution is characterized by five points of support, the mass probability parameters computed by **hsmlogit** take the following expression:  $\pi_j = \frac{e^{p_j}}{1 + \sum_{l=2}^5 e^{p_l}}$ , for  $j = 2, \dots, 5$ , and  $\pi_1 = 1 - \sum_{l=2}^5 \pi_l$ . And for the computation of the standard errors of mass probability parameters, **hsmlogit** also provides to **\_diparm()** command the algebraic expressions of the first order derivatives of each  $\pi_j = \frac{e^{p_j}}{1 + \sum_{l=2}^L e^{p_l}}$ , for each  $j = 1, 2, \dots, P$ , with respect to each  $p_l$ , with  $l = 2, 3, \dots, P$ .

## 4 Command syntax

The **hsmlogit**'s command syntax follows the same design that **hshaz** and **hshaz2**'s. The only difference between **hsmlogit**'s command syntax and **hshaz2**'s is added by **dead(deadvar)** option. Unlike **hshaz2**, the **dead(deadvar)** option of **hsmlogit** command identifies whether the dependent variable typed by the command user in the **dead(deadvar)** option takes one, two or three risks. Therefore, **hsmlogit**, depending on the number of the values taken by the dependent variable, estimates, respectively, a single, a two, or a three competing risks duration model. The rest options of **hsmlogit**'s command syntax are the same that **hshaz2**'s (see [Troncoso-Ponce, 2017]).

The **hsmlogit** command syntax is:

```
hsmlogit varlist [weight] [if exp] [in range] [, id(idvar) dead(deadvar)
    seq(seqvar) spell(spellvar) nmp(#) m2(#) p2(#) m3(#) p3(#)
    m4(#) p4(#) m5(#) p5(#) eform nocons nolog nobeta0 level(#)
    maximizeoptions]
```

<sup>5</sup>As previously mentioned, **hsmlogit** allows for the estimation of a maximum of five points of support (ie. Types of workers) for the identification of the UH distribution.

## 5 Estimation results

This Section shows results from the estimation of three duration models, each of them depends on the number of exits modeled. The first model, presented below in the first estimation output (**Single risk model with UH using hsmlogit**), simply estimates the transition rate out of employment without differentiating the destination state. The second one, shown below in the second estimation output (**Two competing risks model with UH using hsmlogit**), estimates a two risks duration model, by which the two modeled risks are: i) exiting to unemployment; and ii) a job-to-job transition to another employment. Finally, the third model, shown below in the third estimation output (**Three competing risks model with UH using hsmlogit**), allows for distinguishing the type of labor contract of the new employment found in the job-to-job transition. Specifically, the model differs between finding a fixed-term contract, and an open-ended one. Therefore, these three competing risks are: i) exiting to unemployment; ii) finding a fixed-term contract; and iii) finding an open-ended contract.

As mentioned in Section 3, the functional form of the hazard rate estimated in the first model is given by a *Logit* function, whereas the hazard rates of the second and third models are given by *Multinomial Logit* functions with two a three competing risks, respectively. The mentioned three tables with the estimation output show estimation results of fitting multispells duration models with two mass-points of unobserved heterogeneity.<sup>6</sup>

For the three estimated models, the set of covariates included in the specification of the hazard rates controls for the effect of: i) personal characteristics of the employed workers, such as, gender, age (`age16tv`) and squared age (`age16tv2`),<sup>7</sup> nationality,<sup>8</sup> and educational level<sup>9</sup>; ii) business cycle effects, by including the quarterly unemployment rate (`unrate`) and the product of the unemployment rate with the natural logarithm of the current employment spell (`unratexlnjemp`), and its squared (`unratexlnjemp2`); iv) a set of dummy variables that identify the Spanish regions (`andal-rioja`) to capture regional effects. Additionally to the duration dependence specification (using a three order polynomial of the natural logarithm of the duration of current employment spell), three dummy variables are included to identify months 6, 12, 18 and 24. These dummy variables are included to capture exit peaks related to the duration of temporary contracts in the Spanish labor market. Finally, to capture the effect of holding an apprenticeship contract on the employment exit

<sup>6</sup>All estimation results, with and without UH, shown in this article are available to the interested reader upon request.

<sup>7</sup>Age covariates measure the difference between the current age (time-varying age) with respect to the legal working age in the Spanish labor market, 16 years old.

<sup>8</sup>Nationality effect is captured using a dummy variable, called `inmigra`, that takes value one if the employed worker is not Spanish, and zero otherwise.

<sup>9</sup>The effect of educational level is captured by including two dummy variables: `educcompul1` and `educcompul2`. Dummy variable `educcompul1` (`educcompul2`) takes value one whether the worker has a primary (secondary) compulsory education degree, and takes zero otherwise.

rate, the dummy variable ( $cf$ ) takes value one whether the worker is holding an apprenticeship contract, and takes value zero whether the worker has another type of temporary contract different from the apprenticeship one.<sup>10</sup> The regression coefficients not shown in the estimation output tables are omitted due to space reasons, and are available to the interested reader upon request.

---

<sup>10</sup>Hence, the regressions' constant term contains male native employed workers holding a temporary contract in the Spanish regions Madrid and Catalonia, with less than primary compulsory education.

## Single risk model with UH using hsmlogit

```
. hsmlogit_v4 `varsaleE' , id(codind) spell(spell) seq(j) d(exit1) nmp(2) difficult
Discrete time competing risks hazard model without frailty
```

Logistic regression

Number of obs	=	1316611
LR chi2(28)	=	131531.43
Prob > chi2	=	0.0000
Pseudo R2	=	0.0911

Log likelihood = -655754.83

exit1	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
cf	-.8912589	.0084509	-105.46	0.000	-.9078224	-.8746955
lnjump	-1.886904	.018425	-102.41	0.000	-1.923016	-1.850791
lnjump2	.9354424	.0114129	81.96	0.000	.9130735	.9578113
lnjump3	-.1648652	.0024098	-68.41	0.000	-.1695883	-.1601421
age16tv	-.0211124	.0020584	-10.26	0.000	-.0251468	-.017078
age16tv2	.0003848	.0001402	2.75	0.006	.0001101	.0006595
educcompul2	-.0368489	.0044142	-8.35	0.000	-.0455005	-.0281973
manufactory	-.3241511	.0066488	-48.75	0.000	-.3371825	-.3111197
highserv	.1818382	.0088225	20.61	0.000	.1645465	.1991299
lowserv	.1175756	.0060339	19.49	0.000	.1057494	.1294018
unrate	.0002062	.0009034	0.23	0.819	-.0015644	.0019768
unratexlnj-p	.0090781	.0012325	7.37	0.000	.0066624	.0114938
unratexlnj-2	-.0033776	.0004145	-8.15	0.000	-.00419	-.0025653
_cons	.0830947	.0119988	6.93	0.000	.0595775	.1066118

Discrete time competing risks hazard model, with discrete mixture

Log likelihood = -635096

Number of obs	=	1,316,611
---------------	---	-----------

exit1	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
hazard						
cf	-.913632	.0092041	-99.26	0.000	-.9316717	-.8955923
lnjump	-1.437325	.0194947	-73.73	0.000	-1.475534	-1.399116
lnjump2	.7870148	.0118972	66.15	0.000	.7636968	.8103328
lnjump3	-.1451242	.0024777	-58.57	0.000	-.1499804	-.1402681
age16tv	-.0277728	.0023687	-11.73	0.000	-.0324153	-.0231303
age16tv2	.0004827	.0001626	2.97	0.003	.000164	.0008014
educcompul2	-.0291969	.005713	-5.11	0.000	-.0403942	-.0179996
manufactory	-.3281004	.0074911	-43.80	0.000	-.3427826	-.3134182
highserv	.077762	.0100445	7.74	0.000	.0580751	.097449
lowserv	.0705415	.0068966	10.23	0.000	.0570244	.0840586
unrate	.0023788	.0010097	2.36	0.018	.0003999	.0043578
unratexlnjump	.0099583	.0012989	7.67	0.000	.0074125	.0125041
unratexlnjump2	-.0040484	.0004267	-9.49	0.000	-.0048848	-.0032121
_cons	-.4461856	.0145465	-30.67	0.000	-.4746963	-.417675
m2						
_cons	1.588835	.0089093	178.33	0.000	1.571373	1.606297
logitp2						
_cons	-1.742683	.0211864	-82.25	0.000	-1.784207	-1.701158
Prob. Type 1	.8510275	.002686	316.84	0.000	.8456859	.8562156
Prob. Type 2	.1489725	.002686	55.46	0.000	.1437844	.1543141

Note: m1 = 0

## Two competing risks model with UH using hsmlogit

```
. hsmlogit_v4 `varsaleE' , id(codind) spell(spell) seq(j) d(exit2) nmp(2) difficult
Discrete time competing risks hazard model without frailty
```

```
Multinomial logistic regression      Number of obs   =   1316611
                                      LR chi2(56)        =  149758.52
                                      Prob > chi2         =    0.0000
Log likelihood = -861174.84          Pseudo R2        =    0.0800
```

exit2	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
EU						
cf	-.7949866	.0109422	-72.65	0.000	-.816433	-.7735403
lnjemp	-.9706289	.0245682	-39.51	0.000	-1.018782	-.922476
lnjemp2	.5054906	.0152706	33.10	0.000	.4755607	.5354205
lnjemp3	-.1151439	.0033619	-34.25	0.000	-.1217332	-.1085547
age16tv	-.0874133	.0026712	-32.72	0.000	-.0926488	-.0821778
age16tv2	.0039427	.0001826	21.60	0.000	.0035849	.0043005
educcompul2	-.1047359	.0060011	-17.45	0.000	-.1164979	-.0929739
manufactory	-.1986758	.0088302	-22.50	0.000	-.2159827	-.1813689
highserv	.2298115	.0118434	19.40	0.000	.2065989	.2530242
lowserv	.2196708	.0079353	27.68	0.000	.204118	.2352237
unrate	.0225876	.0012259	18.43	0.000	.0201849	.0249904
unratexlnj-p	-.0009619	.0016357	-0.59	0.557	-.0041678	.0022441
unratexlnj-2	.0002686	.0005461	0.49	0.623	-.0008018	.0013389
_cons	-1.011903	.0163099	-62.04	0.000	-1.04387	-.9799364
EE						
cf	-1.011942	.0122039	-82.92	0.000	-1.035861	-.9880226
lnjemp	-2.53216	.0241991	-104.64	0.000	-2.57959	-2.484731
lnjemp2	1.255083	.0153547	81.74	0.000	1.224988	1.285178
lnjemp3	-.2022639	.0031734	-63.74	0.000	-.2084836	-.1960442
age16tv	.043099	.0027103	15.90	0.000	.037787	.0484111
age16tv2	-.0031777	.000185	-17.18	0.000	-.0035402	-.0028152
educcompul2	.0220555	.005569	3.96	0.000	.0111405	.0329705
manufactory	-.4401769	.0088543	-49.71	0.000	-.4575309	-.4228228
highserv	.1415707	.0109013	12.99	0.000	.1202046	.1629368
lowserv	.0282416	.0076963	3.67	0.000	.0131571	.0433262
unrate	-.0165131	.0010909	-15.14	0.000	-.0186511	-.0143751
unratexlnj-p	.0118484	.0016218	7.31	0.000	.0086697	.0150271
unratexlnj-2	-.0052452	.0005698	-9.21	0.000	-.006362	-.0041284
_cons	-.3861184	.0145517	-26.53	0.000	-.4146391	-.3575977

(exit2==no exit is the base outcome)

Discrete time competing risks hazard model, with discrete mixture

Log likelihood = -840603.57

Number of obs = 1,316,611

exit2	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
hazard1						
cf	-.8244588	.0115585	-71.33	0.000	-.8471131	-.8018045
lnjemp	-.5313624	.0253333	-20.97	0.000	-.5810148	-.4817099
lnjemp2	.3610352	.01561	23.13	0.000	.3304402	.3916302
lnjemp3	-.0957306	.0034033	-28.13	0.000	-.102401	-.0890603
age16tv	-.0934056	.0028976	-32.24	0.000	-.0990848	-.0877263
age16tv2	.0039949	.0001992	20.06	0.000	.0036046	.0043852
educcompul2	-.0968852	.0069518	-13.94	0.000	-.1105105	-.08326
manufactory	-.211835	.0094664	-22.38	0.000	-.2303887	-.1932812
highserv	.1350418	.0127147	10.62	0.000	.1101215	.1599621
lowserv	.1745412	.0085626	20.38	0.000	.1577589	.1913235
unrate	.024601	.0013045	18.86	0.000	.0220443	.0271577
unratexlnjemp	.0004602	.0016838	0.27	0.785	-.0028401	.0037605
unratexlnjemp2	-.0005451	.0005548	-0.98	0.326	-.0016324	.0005422
_cons	-1.540267	.0182031	-84.62	0.000	-1.575944	-1.504589
hazard2						
cf	-1.02528	.0127572	-80.37	0.000	-1.050284	-1.000277
lnjemp	-2.078786	.0250507	-82.98	0.000	-2.127885	-2.029688
lnjemp2	1.107497	.015722	70.44	0.000	1.076683	1.138312
lnjemp3	-.1831092	.0032193	-56.88	0.000	-.1894189	-.1767995
age16tv	.038407	.0029839	12.87	0.000	.0325588	.0442553
age16tv2	-.0031854	.000205	-15.54	0.000	-.0035871	-.0027836
educcompul2	.0336721	.0067271	5.01	0.000	.0204872	.0468571
manufactory	-.4384897	.0095412	-45.96	0.000	-.45719	-.4197893
highserv	.0295985	.0119948	2.47	0.014	.006089	.0531079
lowserv	-.0230854	.0084529	-2.73	0.006	-.0396528	-.0065181
unrate	-.0144479	.0011802	-12.24	0.000	-.0167611	-.0121346
unratexlnjemp	.0123513	.0016725	7.38	0.000	.0090732	.0156294
unratexlnjemp2	-.0058043	.0005786	-10.03	0.000	-.0069384	-.0046703
_cons	-.9270718	.0167796	-55.25	0.000	-.9599592	-.8941844
m2						
_cons	1.584087	.0088521	178.95	0.000	1.566737	1.601437
logitp2						
_cons	-1.730903	.0210766	-82.12	0.000	-1.772212	-1.689593
Prob. Type 1	.8495279	.0026942	315.31	0.000	.8441707	.8547326
Prob. Type 2	.1504721	.0026942	55.85	0.000	.1452674	.1558293

Note: m1 = 0



### Three competing risks with UH using hsmlogit

```
. hsmlogit_v4 `varsaleE' , id(codind) spell(spell) seq(j) d(exit3) nmp(2) difficult
Discrete time competing risks hazard model without frailty
```

Multinomial logistic regression

Number of obs	=	1316611
LR chi2(84)	=	160654.70
Prob > chi2	=	0.0000
Pseudo R2	=	0.0836

Log likelihood = -881090.49

exit3	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
EU						
cf	-.7986611	.0109406	-73.00	0.000	-.8201043	-.7772179
lnjemp	-.9709934	.0245729	-39.51	0.000	-1.019155	-.9228314
lnjemp2	.5069499	.0152775	33.18	0.000	.4770066	.5368933
lnjemp3	-.1157149	.0033656	-34.38	0.000	-.1223114	-.1091185
age16tv	-.0881035	.0026718	-32.98	0.000	-.0933401	-.0828668
age16tv2	.0039923	.0001826	21.86	0.000	.0036344	.0043502
educcompul2	-.105202	.0060017	-17.53	0.000	-.1169652	-.0934388
manufactory	-.2006661	.0088305	-22.72	0.000	-.2179736	-.1833587
highserv	.2296968	.0118453	19.39	0.000	.2064804	.2529132
lowserv	.2180955	.0079363	27.48	0.000	.2025406	.2336503
unrate	.022578	.0012261	18.41	0.000	.0201749	.0249811
unratexlnj-p	-.0011119	.0016364	-0.68	0.497	-.0043191	.0020953
_cons	-1.009577	.0163131	-61.89	0.000	-1.04155	-.9776039
ET						
cf	-1.214204	.0134696	-90.14	0.000	-1.240604	-1.187804
lnjemp	-2.566013	.0252271	-101.72	0.000	-2.615457	-2.516569
lnjemp2	1.310668	.0165322	79.28	0.000	1.278266	1.343071
lnjemp3	-.2269325	.0035308	-64.27	0.000	-.2338528	-.2200123
age16tv	.0348604	.0027452	12.70	0.000	.0294799	.0402409
age16tv2	-.002567	.0001871	-13.72	0.000	-.0029338	-.0022002
educcompul2	.0164089	.0056687	2.89	0.004	.0052985	.0275193
manufactory	-.4715925	.0091037	-51.80	0.000	-.4894354	-.4537496
highserv	.1380367	.0110533	12.49	0.000	.1163726	.1597009
lowserv	.0082154	.0078295	1.05	0.294	-.0071302	.023561
unrate	-.0164156	.001101	-14.91	0.000	-.0185736	-.0142576
unratexlnj-p	.0097402	.0016896	5.76	0.000	.0064288	.0130517
_cons	-.3760623	.0146951	-25.59	0.000	-.4048642	-.3472603
EP						
cf	.7109967	.035604	19.97	0.000	.6412142	.7807791
lnjemp	-2.083359	.1148221	-18.14	0.000	-2.308406	-1.858312
lnjemp2	1.077243	.056833	18.95	0.000	.9658519	1.188633
lnjemp3	-.1234547	.0099939	-12.35	0.000	-.1430423	-.1038671
age16tv	.2424454	.0156192	15.52	0.000	.2118325	.2730584
age16tv2	-.0191269	.0011496	-16.64	0.000	-.0213801	-.0168737
educcompul2	.151788	.0268159	5.66	0.000	.0992299	.2043462
manufactory	.149392	.0355004	4.21	0.000	.0798125	.2189714
highserv	.2252933	.057029	3.95	0.000	.1135184	.3370681
lowserv	.5072968	.0360027	14.09	0.000	.4367328	.5778609
unrate	-.0802567	.0089966	-8.92	0.000	-.0978896	-.0626237
unratexlnj-p	.0610449	.0089799	6.80	0.000	.0434446	.0786453
_cons	-4.808289	.1012222	-47.50	0.000	-5.006681	-4.609897

(exit3==no exit is the base outcome)

Discrete time competing risks hazard model, with discrete mixture

Log likelihood = -860570.7

Number of obs = 1,316,611

exit3	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
hazard1						
cf	-.8317648	.0115702	-71.89	0.000	-.8544421	-.8090876
lnjemp	-.5323972	.0253395	-21.01	0.000	-.5820617	-.4827327
lnjemp2	.3637062	.0156171	23.29	0.000	.3330973	.3943152
lnjemp3	-.0965266	.0034067	-28.33	0.000	-.1032035	-.0898496
age16tv	-.0946958	.0028946	-32.71	0.000	-.1003691	-.0890224
age16tv2	.0041084	.0001986	20.68	0.000	.003719	.0044977
educcompul2	-.098385	.0069564	-14.14	0.000	-.1120192	-.0847508
manufactory	-.2151223	.0094659	-22.73	0.000	-.2336751	-.1965695
highserv	.1342664	.0127151	10.56	0.000	.1093452	.1591875
lowserv	.1714203	.0085664	20.01	0.000	.1546305	.1882101
unrate	.0246167	.0013044	18.87	0.000	.02206	.0271733
unratexlnjemp	.0002379	.0016847	0.14	0.888	-.003064	.0035398
_cons	-1.537957	.0182095	-84.46	0.000	-1.573647	-1.502267
hazard2						
cf	-1.230282	.0139993	-87.88	0.000	-1.25772	-1.202843
lnjemp	-2.109857	.0260808	-80.90	0.000	-2.160975	-2.05874
lnjemp2	1.161366	.0169049	68.70	0.000	1.128233	1.194499
lnjemp3	-.2072636	.0035771	-57.94	0.000	-.2142745	-.2002527
age16tv	.0294043	.0030183	9.74	0.000	.0234885	.0353201
age16tv2	-.002497	.000207	-12.06	0.000	-.0029027	-.0020913
educcompul2	.0270498	.0068402	3.95	0.000	.0136433	.0404563
manufactory	-.4710222	.0097941	-48.09	0.000	-.4902183	-.4518261
highserv	.023543	.0121529	1.94	0.053	-.0002762	.0473622
lowserv	-.0459738	.008593	-5.35	0.000	-.0628157	-.0291318
unrate	-.0143704	.0011898	-12.08	0.000	-.0167023	-.0120384
unratexlnjemp	.0101926	.0017401	5.86	0.000	.006782	.0136031
_cons	-.9171136	.0169239	-54.19	0.000	-.9502838	-.8839434
hazard3						
cf	.6758236	.0358259	18.86	0.000	.6056062	.7460411
lnjemp	-1.679314	.1147489	-14.63	0.000	-1.904217	-1.45441
lnjemp2	.9649278	.0568303	16.98	0.000	.8535425	1.076313
lnjemp3	-.1111428	.0099834	-11.13	0.000	-.13071	-.0915757
age16tv	.2357879	.0156333	15.08	0.000	.2051473	.2664286
age16tv2	-.0190263	.0011489	-16.56	0.000	-.0212781	-.0167745
educcompul2	.1549394	.0269799	5.74	0.000	.1020597	.2078191
manufactory	.134427	.035673	3.77	0.000	.0645092	.2043449
highserv	.1660599	.0571961	2.90	0.004	.0539576	.2781622
lowserv	.4844141	.0361505	13.40	0.000	.4135604	.5552679
unrate	-.0763687	.0089703	-8.51	0.000	-.0939502	-.0587873
unratexlnjemp	.0608414	.0089472	6.80	0.000	.0433051	.0783776
_cons	-5.338236	.101233	-52.73	0.000	-5.536649	-5.139823
m2						
_cons	1.583831	.0088365	179.24	0.000	1.566512	1.60115
logitp2						
_cons	-1.725456	.0210304	-82.05	0.000	-1.766674	-1.684237
Prob. Type 1	.8488302	.0026986	314.55	0.000	.8434647	.8540436
Prob. Type 2	.1511698	.0026986	56.02	0.000	.1459564	.1565353

Note:  $m1 = 0$

The estimation exercise shown in this Section is addressed only to highlight the importance of allowing for modelling more than one single risk in a duration model that also takes into account the presence of UH. For that reason, analogously to [Troncoso-Ponce, 2017], the main purpose of these regressions is not intended to address a rigorous regression analysis to properly estimate the effect of a set of covariates on the probability of exiting out of employment. Therefore, in this Section, comments on detailed estimation results will be focused mainly on the impact of holding an apprenticeship contract (captured by the covariate *cf* in the three estimation outputs presented above) when we allow for modelling more than one single risk.

The single risk duration model estimates a statistically significant negative effect ( $-0.9136$ ) of holding an apprenticeship contract on the probability of exiting out from the employment state, which may suggest that apprenticeship contracts last longer (ie. seem to be more stable) than regular fixed-term contracts. And when we allow for modelling two competing risks (exiting to unemployment, or a job-to-job transition to another job), the effect of apprenticeship contracts remain negative and statistically significant on both the two risks modeled: exiting to exit to unemployment ( $-0.8244$ ), and a direct transition to another job ( $-1.0252$ ).

However, interestingly, the estimated effect of apprenticeship contracts turns positive when we allow for modelling the job-to-job transition separately in two different, and mutually exclusive, destinations: i) a direct transition to a fixed-term contract; and i) a direct transition to an open-ended contract. As the third estimation output shows, apprenticeship contracts increase the probability of experiencing a job-to-job transition towards an open-ended contract ( $0.7109$ ). The main reason of observing this positive effect is the role played by public financial incentives addressed to the conversion of apprenticeship contracts into open-ended ones. Apprenticeship contracts in Spain benefit from public subsidies for the conversion into open-ended contracts. These subsidies mainly consist of a significant reduction in Social Security contributions paid by the employer during a maximum period of three years, from the starting date of conversion of the apprenticeship contract into an open-ended one. The main goal of these financial incentives is to favour employment stability, and to foster the accumulation of employment experience of apprentices by allowing them to put in practice the work-specific skills acquired during the apprenticeship period. Thus, the positive coefficient found ( $0.7109$ ) may be capturing the effect of these public financial incentives provided by Spanish policy makers addressed to the conversion of apprenticeship contracts into open-ended ones.

An exhaustive analysis of the apprenticeship contracts in the Spanish labor market is presented in [Troncoso-Ponce, 2016] and [Jansen and Troncoso-Ponce, 2017]. The first one estimates a multispell and multistate competing risks duration

Table 1: Interpretation of UH coefficients (three competing risks model)

	Prob.	Emp. to Unemp.	Emp. to Fixed-term	Emp. to Open-ended
Type I	84.88%	-1.537957	-0.9171136	-5.338236
Type II	15.12%	0.045874 <sup>a</sup>	0.6667174 <sup>b</sup>	-3.754405 <sup>c</sup>
<sup>a</sup> (= -1.537957 + 1.583831) <sup>b</sup> (= -0.9171136 + 1.583831) <sup>c</sup> (= -5.338236 + 1.583831)				

model with UH specific to both each state and to each destination state, as well as a selection equation that estimates the transition rates to the entry into the labor market holding three different types of labor contract: an apprenticeship contract, a fixed-term contract, and an open-ended contract. The second, and more recent, work also estimates a multispell and multistate competing risks duration model with UH, but the selection equation consists of an initial conditions equation, rather than a transition rate equation, that controls for the effect of a set of observable covariates on the probability of having an apprenticeship contract just in the first employment spell of the individual's working life. Moreover, the empirical strategy followed in this work allows us to disentangle two types of effect: an instant effect, and a subsequent effect of apprenticeship contracts on the employment and unemployment transition rates.

## 5.1 Some insights on the interpretation of UH coefficients

Regarding the estimation and interpretation of UH coefficients, as we assume that  $\eta_1$  is set to zero,<sup>11</sup> the estimated regression's constant terms (-1.537957, -0.9171136 and -5.338236, for the exit to unemployment, to a fixed-term, and to an open-ended contract, respectively) capture the UH component specific to Type I workers, whereas  $\eta_2$  captures the unobserved differential effect of Type II workers with respect to Type I workers. Therefore, the estimated value of UH component specific to Type II workers are the result of the sum of the regression's constant terms and the estimated coefficient value of  $\eta_2$ .

Table 1 shows the estimated coefficients of the UH components of Type I and Type II workers from the estimation results of the three competing risks model. The estimation of the non-parametric UH distribution, characterized by the presence of two types of workers (two points of support), captures Type I and Type II workers who represent, respectively, 84.88% and 15.12% of the estimation sample. As Table 1 shows, Type II workers have unobserved characteristics that positively correlate to the employment hazard rates, which implies that Type II workers face employment transition rates (towards all the three

<sup>11</sup>It explains the footnote shown at the estimation output tables with the message "Note: m1 = 0", where m1 denotes UH component given by  $\eta_1$ . See also [Troncoso-Ponce, 2017] and `hshaz` command's official Stata helpfile.

modeled risks) higher than Type I workers’.

In conclusion, the estimation of not only a single or a two competing risks, but a three competing risks duration model has allowed for capturing a positive and statistically significant effect of apprenticeship contracts on the probability of transiting directly (via job-to-job) to an open-ended contract, that otherwise would have remained hidden to the empirical researcher if only one risk, or even two, would have been estimated. Furthermore, given the relevance of UH, and its non parametric identification, in discrete time duration models with multispell observations (see, for example, [Gaure, Roed and Zhang, 2007] and [Abbring and Van den Berg, 2004]), the new Stata command `hsmlogit` takes especial relevance, as allows for the estimation of discrete time competing risks duration models with UH.

## 6 The advantages of using `ml d2` method

As mentioned in Section 1, `hsmlogit` provides the algebraic expressions of both the gradient vector and Hessian matrix, allowing for using `d2 ml` method to achieve the model convergence. An important advantage of programming the Hessian matrix is that allows applied researchers to deal with large longitudinal microdata sets (see for example [Troncoso-Ponce, 2017]). To show the savings in estimation time, this Section presents time required to estimate multispell both single and competing risks duration models (with 2, 3 and 4 mass-points) using `d0`, `d1` and `d2 ml` methods,<sup>12</sup>. Comments in this Section will be focused only on the comparison between `d1` and `d2 ml` methods. The comparison between `d0` and `d2 ml` does reinforce the same conclusions obtained below.

Table 2 reports time spent<sup>13</sup> by each of the three `ml` methods in achieving the models’ convergence.<sup>14</sup> Results from Table 2 highlight two relevant differences between `d1` and `d2 ml` methods: Firstly, `d2` method significantly reduces time required to achieve the all models convergence. Differences in time required seem to be less evident in the estimation of single risk models: for instance, for fitting the two mass-points model, `d2 (d1)` method needs 46 seconds (6.28 minutes). However it becomes more important as both the number of risks and the number of mass-points increase: for fitting the three competing risks model with four mass-points, `d2` method only requires 8.02 minutes, whereas `d1` method needs 1.58 hours. On its part, `d0` method not even achieve the model convergence:

<sup>12</sup>The all estimations, whose time required are shown in Table 2, include a set of twenty eight covariates that, as results shown in Section 5, control for duration dependence, personal characteristics, type of labor contract, regional effects and economic cycle. The detailed estimation results are available upon request to the interested reader.

<sup>13</sup>I work with Stata 14.0 MP - Parallel edition 64 bits. The machine employed to obtain estimation results incorporates an Intel(R) Core(TM) i7-6700HQ CPU at 2.60 GHz, and 12 Gb RAM memory. The operating system is Windows 10 Home.

<sup>14</sup>In this sample composed of youth Spanish employees, the estimation of five points of support for the identification of the non-parametric unobserved heterogeneity distribution is not possible, neither fitting single risk models, nor competing risks models.

Table 2: Time required for the estimation of multispell competing risks duration models (Sample size: 1,316,611 observations)

	Time (hh:mm:ss)				
	d0 method	d1 method	d2 method	Diff.=d1-d2	Diff.=d0-d2
Single risk					
Two mass-points	1:37:59	0:06:28	0:00:46	0:05:42	1:32:17
Three mass-points	3:00:16	0:12:48	0:02:03	0:10:45	2:49:31
Four mass-points	18:02:16	0:21:01	0:06:52	0:14:09	17:48:07
Two risks					
Two mass-points	3:59:59	0:20:10	0:01:48	0:18:22	3:41:37
Three mass-points	7:58:44	0:37:28	0:04:00	0:33:28	7:25:16
Four mass-points	7:03:51	0:58:28	0:07:06	0:51:22	6:12:29
Three risks					
Two mass-points	3:45:53	0:40:09	0:03:19	0:36:50	3:09:03
Three mass-points	-	1:13:42	0:05:43	1:07:59	-
Four mass-points	-	1:58:29	0:08:02	1:50:27	-

after eleventh iteration, it gets into a backed up loop. Secondly, unlike **d2** method, time required by **d1** method to achieve the model convergence strongly depends both on the number of exits modeled, and on the number of points of support for the identification of the UH. Table 2 shows that, using **d2** (**d1**) method, the difference in time spent between the less time-demanding model (the single risk model with two mass-points) and the most time-demanding model (the three competing risks with four mass-points) reaches 7.16 minutes (1.52 hours).

## 7 Concluding remarks

This article presents **hsmlogit**, a new Stata command that estimates multispells discrete time competing risks duration models with unobserved heterogeneity. **hsmlogit** allows for the estimation of one, two and up to three competing risks, as well as a maximum of five points of support for the identification of the non-parametric unobserved heterogeneity distribution [Heckman and Singer, 1984]. The relevance of modelling more than one risk has been highlighted by estimating the effect of apprenticeship contracts on a sample composed of low educated young workers in the Spanish labor market for the period 2000-2014. Thus, the estimation of a three competing risks duration model has been the only way to find out the potential effect of public financial incentives for the conversion of apprenticeship contracts into open-ended ones on the direct (via job-to-job) transition rates towards an open-ended contract. Moreover, since **hsmlogit** allows for the estimation of non-parametric UH distribution ([Heckman and Singer, 1984]),

our results capture the presence of two types of workers with different values of unobserved characteristics that affect the estimated hazard rates.

Finally, `hsmlogit` provides the algebraic expressions of both the gradient vector and the Hessian matrix, which significantly reduces time required to achieve the model convergence, and also improves the standard errors' accuracy of the estimated coefficients. The possibility of estimating competing risks duration models with the presence of UH, along with time savings provided by the use of `d2 ml` method may allow the applied researchers to easily and properly exploit the richness and complexity of large longitudinal microdata sets.

## References

- [Abbring and Van den Berg, 2004] ABBRING, JAAP H. AND VAN DEN BERG, GERARD J., *Analysing the effect of dynamically assigned treatments using duration models, binary treatment models, and panel data models*, Empirical Economics, January 2004, Volume 29, Issue 1, pp. 5-20.
- [Allison, 1982] ALLISON, PAUL D., *Discrete-Time Methods for the Analysis of Event Histories*, Sociological Methodology, Vol. 13 (1982), pp. 61-98.
- [Arranz, García-Serrano and Hernanz, 2013] ARRANZ, J.M., GARCÍA-SERRANO, C. AND HERNANZ, V., 2013, *How do we pursue ?labormetrics?? An application using the MCVL*, Estadística Española, Vol. 55 (2013), No. 181, pp. 231-254.
- [Arranz and García-Serrano, 2011] ARRANZ, J.M. AND GARCÍA-SERRANO, C., 2011, *Are the MCVL tax data useful? Ideas for mining*, Hacienda Pública Española, Vol. 199(4), pp. 151-186.
- [García-Pérez, 2008] GARCÍA-PÉREZ, J.I., 2008 *La Muestra Continua de Vidas Laborales: Una guía de uso para el análisis de transiciones*, Revista de Economía Aplicada, N. E-1, Vol. XVI, pp. 5-28.
- [Gaure, Roed and Zhang, 2007] GAURE, S., ROED, K. AND ZHANG, T., 2007 *Time and causality: A Monte Carlo assessment of the timing-of-events approach*, Journal of Econometrics, Volume

- 141, Issue 2, December 2007, Pages 1159-1195.
- [Gould et al., 2010] GOULD, W., PITBLADO, J. AND POI, B., *Maximum Likelihood Estimation with Stata*, Fourth Edition, Stata Press, 2010.
- [Heckman and Singer, 1984] HECKMAN, J. J. AND SINGER, B., *A Method for Minimizing the Impact of the Distributional Assumptions in Econometric Models for Duration Data*, *Econometrica*, Vol. 52, pp. 271-320.
- [Jansen and Troncoso-Ponce, 2017] JANSEN, M. AND TRONCOSO-PONCE, D., *The impact of apprenticeship contracts on the labour market insertion of youth in Spain*, Fedea working paper (forthcoming).
- [Jenkins, 1995] JENKINS, S., 1995 *Easy Estimation Methods for Discrete-Time Duration Models*, *Oxford Bulletin of Economics and Statistics*, Vol. 57, 1, 1995.
- [Lancaster, 1992] LANCASTER, TONY, *The Econometric Analysis of Transition Data*, First Edition, Cambridge University Press, 1992.
- [Lapuerta, 2010] LAPUERTA, I. (2010) *Claves para el trabajo con la Muestra Continua de Vidas Laborales*, DemoSoc working paper (2010-37), Universitat Pompeu Fabra
- [Troncoso-Ponce, 2016] TRONCOSO-PONCE, D., 2016 *An empirical analysis of some public policies applied to the Spanish labour market*, PhD. Thesis, Chapter 3.
- [Troncoso-Ponce, 2017] TRONCOSO-PONCE, D., 2017 *Faster estimation of discrete time duration models using hshaz2*, Manuscript. Available at: <https://ideas.repec.org/p/pab/wpaper/17.05.html>
- [Van Den Berg, 2001] VAN DEN BERG, GERARD. J, 2001 *Duration models: specification, identification, and multiple durations*, *Handbook of Econometrics*, Elsevier, Vol. 5, 2001, pp. 3381-3460.