

The prediction of time to default for personal loans using mixture cure models: including macro-economic factors.

L. Dirick¹ T. Bellotti²
G. Claeskens¹ B. Baesens^{1,3}

¹ Faculty of Economics and Business, KU Leuven, Belgium

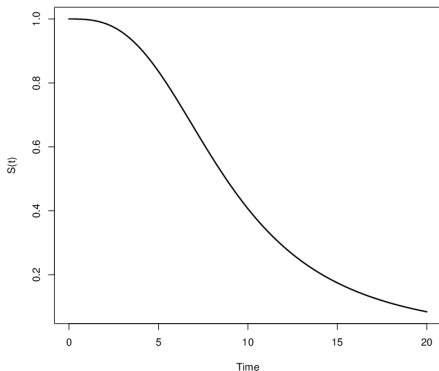
² Department of Mathematics, Imperial College, London, UK

³ School of Management, University of Southampton, UK

Credit scoring & credit control XIV conference, Edinburgh

Why survival analysis?

- ▶ Basel II accord: more accurate credit risk calculations needed.
- ▶ Survival analysis: event of interest = default.
 - compute profitability over a customer's lifetime.
 - intuitive: possibility to model time aspect.



Survival analysis: issues

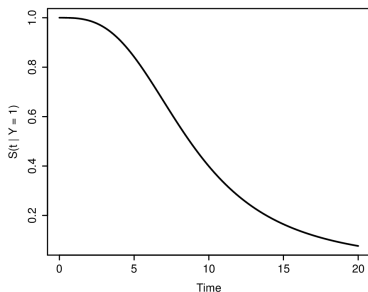
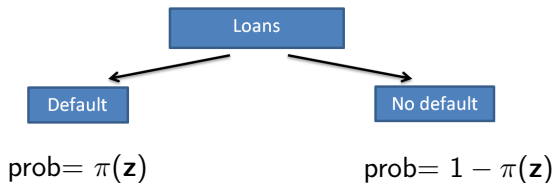
- Idea of survival in credit risk context \neq medical survival.
- “The event might not be observed, but will take place eventually”.
- For credit risk events such as default, this is not the case.

⇒ Solution: a mixture cure model to model *long-term survivors*.

The classical mixture cure model

Given covariate vector \mathbf{z} , susceptibility indicator Y

$$S(t | \mathbf{z}) = \pi(\mathbf{z})S(t | Y = 1, \mathbf{z}) + 1 - \pi(\mathbf{z})$$



in a **classical** mixture cure model (fixed loan term data, here 36 months).

subject	time	δ	z_1	...	z_m
1	22	1	"married"	...	3
2	36	0	"divorced"	...	2
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
n	7	0	"widowed"	...	6



With **time-varying covariates (TVCs)**, in our paper macro-economic factors

subject	time	δ	z_1	...	z_m	$x_1(t=1)$...	$x_1(t=36)$...	$x_l(t=1)$...	$x_l(t=36)$
1	22	1	"married"	...	3	0.3	...	0.5	...	3	...	1
2	36	0	"divorced"	...	2	0.7	...	0.3	...	2	...	2
\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\ddots	\vdots	\ddots	\vdots	\ddots	\vdots
n	7	0	"widowed"	...	6	0.2	...	0.9	...	4	...	3

δ : 1 if default, 0 if **no default** OR **not yet observed**.

subject	interval	time	λ	δ	z_1	...	z_m	$x_1(t)$...	$x_l(t)$
1	[0, 1]	22	0	1	"married"	...	3	0.3	...	3
1	[1, 2]	22	0	1	"married"	...	3	0.6	...	5
1	\vdots	22	0	1	"married"	\ddots	\vdots	\vdots	\ddots	\vdots
1	[21, 22]	22	1	1	"married"	...	3	0.5	...	1
2	[0, 1]	36	0	0	"divorced"	...	2	0.7	...	2
2	[1, 2]	36	0	0	"divorced"	...	2	0.1	...	3
2	\vdots	36	0	0	"divorced"	\ddots	\vdots	\vdots	\ddots	\vdots
2	[35, 36]	36	0	0	"divorced"	...	2	0.3	...	2
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots	\ddots	\vdots
n	[0, 1]	4	0	0	"widowed"	...	6	0.2	...	4
n	[1, 2]	4	0	0	"widowed"	...	6	0.1	...	1
n	\vdots	4	0	0	"widowed"	\ddots	\vdots	\vdots	\ddots	\vdots
n	[6, 7]	4	0	0	"widowed"	...	6	0.9	...	3

Look at each subject on a monthly basis

λ : time-specific censoring indicator when TVCs in a mixture cure model

Model:

$$S(t \mid \mathbf{z}, \mathbf{x}(t)) = \pi(\mathbf{z})S(t \mid Y = 1, \mathbf{z}, \mathbf{x}(t)) + 1 - \pi(\mathbf{z})$$

$\pi(\mathbf{z}) = P(Y = 1 \mid \mathbf{z})$ = the *Incidence model component*, gives proportion of accounts susceptible to default, given $\mathbf{z} = (z_1, \dots, z_m)$.

$$\log\left(\frac{\pi(\mathbf{z})}{1 - \pi(\mathbf{z})}\right) = \mathbf{z}'\mathbf{b} \Leftrightarrow \pi(\mathbf{z}') = \frac{\exp(\mathbf{z}'\mathbf{b})}{1 + \exp(\mathbf{z}'\mathbf{b})}$$

$S(t \mid Y = 1, \mathbf{z}, \mathbf{x}(t))$ = the *Latency model component*, a conditional survival function, given $\mathbf{z} = (z_1, \dots, z_m)$ **AND** $\mathbf{x}(t) = (x_1(t), \dots, x_l(t))$

$$S(t \mid Y = 1, \mathbf{z}, \mathbf{x}(t)) = \exp\left(-\exp(\beta'\mathbf{z} + \beta'_t\mathbf{x}(t)) \int_0^t h_0(u \mid Y = 1)du\right),$$

We are talking about susceptibility, however, susceptibility is not fully observed.

When $\delta = 0$: will default eventually occur or not?

$$Y = \begin{cases} 0 & \text{when an account is non-susceptible to default} \\ 1 & \text{when the account is susceptible to default} \end{cases}$$

Three possible states:

$$\begin{cases} \delta = 1 \text{ and } Y = 1 \\ \delta = 0 \text{ and } Y = 1 \\ \delta = 0 \text{ and } Y = 0 \end{cases}$$

δ fully observed, but Y only observed when $\delta = 1$.

Incomplete information: iterative EM-algorithm.

The likelihood function

For each time interval j from observation i , following information:
 $\mathbf{O}_{ij} = (\mathbf{x}_i(t=j), \mathbf{z}_i, t_i, \delta_i, \lambda_{i,j})$. The complete likelihood is given by:

$$L(\mathbf{b}, \boldsymbol{\beta}) = \prod_{i=1}^n \pi(\mathbf{z}_i)^{Y_i} (1 - \pi(\mathbf{z}_i))^{(1-Y_i)} \\ \times \prod_{j=1}^{k_i} h(t_i | Y_i = 1, \mathbf{x}_i)^{\lambda_{i,j} Y_i} S(t_i | Y_i = 1, \mathbf{x}_i)^{Y_i}$$

then the complete-data log $L_c(\mathbf{b}, \boldsymbol{\beta}, \boldsymbol{\beta}_t | \mathbf{z})$ is the sum of

$$\log L_{inc}(\mathbf{b} | \mathbf{z}, Y) = \sum_{i=1}^n (1 - Y_i)(1 - \pi(\mathbf{z}_i)) + Y_i \pi(\mathbf{z}_i)$$

$$\log L_{lat}(\boldsymbol{\beta}, \boldsymbol{\beta}_t | \mathbf{z}, \mathbf{x}(t), Y) = \sum_{i=1}^n \sum_{j=1}^{k_i} Y_i \lambda_{i,j} \log h(t=j | Y_i=1, \mathbf{z}_i, \mathbf{x}_i(t=j)) \\ + Y_i \log S(t=j | Y_i=1, \mathbf{z}_i, \mathbf{x}_i(t=j)).$$

The EM-algorithm

The p-th Expectation Step

Conditional probability of i -th observation being susceptible

Survival contributions of subject i : TVC-values at censoring time (t_j)

$$w_i^{(p)} = E(Y_i | \Theta^{(p)}, \mathbf{O}) = \begin{cases} \frac{\pi(\mathbf{z}_i)S(t_i | Y_i = 1, \mathbf{z}_i, \mathbf{x}_i(t_i))}{\pi(\mathbf{z}_i)S(t_i | Y_i = 1, \mathbf{z}_i, \mathbf{x}_i(t_i)) + (1 - \pi(\mathbf{z}_i))} & \text{if } \delta_i = 0 \\ 1 & \text{if } \delta_i = 1. \end{cases}$$

The p-th Maximization Step

Maximizing $Q(\Theta^{(p+1)} | \Theta^{(p)}) = E[\log L_c(\mathbf{b}, \boldsymbol{\beta}, \boldsymbol{\beta}_t | \mathbf{z})]$,

parameter update $\hat{\Theta}^{(p+1)} = (\hat{\mathbf{b}}^{(p+1)}, \hat{\boldsymbol{\beta}}^{(p+1)}, \hat{\boldsymbol{\beta}}_t^{(p+1)}, \hat{\mathbf{S}}_0^{(p+1)})$.

Denote $R(t_j)$ the individuals at risk in interval j , then

$$\hat{\mathbf{S}}_0^{(p+1)}(t) = \exp\left(-\sum_{j:t_j \leq t} \frac{\sum_{i \in R(t_j)} \lambda_{i,j}}{\sum_{i \in R(t_j)} w_i^{(p)} \exp(\boldsymbol{\beta}'^{(p)} \mathbf{z}_i + \boldsymbol{\beta}_t'^{(p)} \mathbf{x}_i(t=j))}\right).$$

repeat until convergence.

- ▶ 100 replications and $n = 1000$.
- ▶ Event times exponentially distributed, rate= 0.7.
- ▶ Censoring times exponentially distributed, rate= 0.1/0.2 (Setting 1/2).
- ▶ Simulating TVCs: used method proposed by Austin (2012).
- ▶ Resulting censoring and unsusceptibility %:
 - 32.9 and 22.7 for Setting 1
 - 86.1 and 80.7 for Setting 2

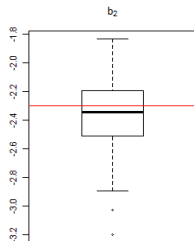
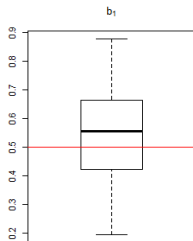
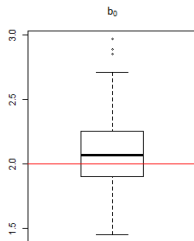
Covariate distributions

variable	z_1	z_2	$x_1(t)$	$x_2(t)$
Distr.	$N(1.5, 0.6)$	$\text{bin}(1, 0.5)$	$N(2, 0.5)$	$N(0.8, 0.5)$

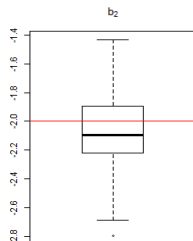
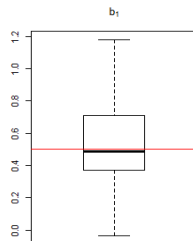
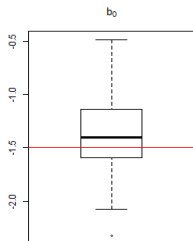
Parameter values of the true model

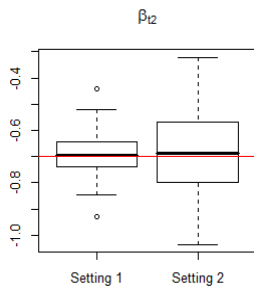
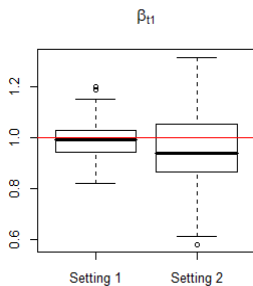
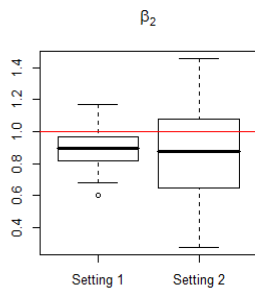
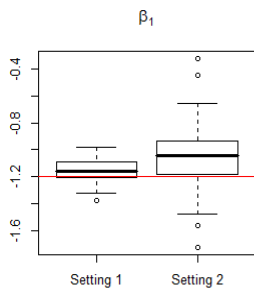
parameter	b_0	b_1	b_2	β_1	β_2	β_{t1}	β_{t2}
Setting 1	2	0.5	-2.3	-1.2	1	1	-0.7
Setting 2	-1.5	0.5	2	-1.2	1	1	-0.7

► Setting 1



► Setting 2



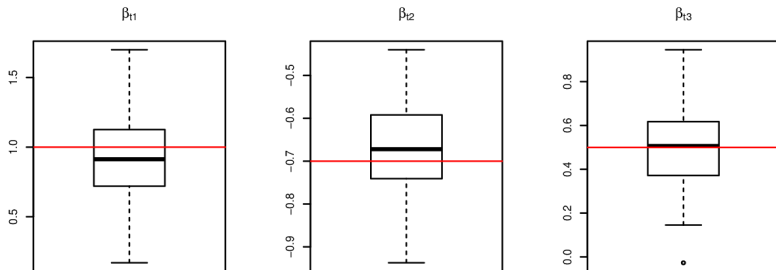


Test for correlated TVCs. Setting 3 = Setting 1, but ...

$$x(t) = \begin{pmatrix} x_1(t) \\ x_2(t) \\ x_3(t) \end{pmatrix}; \quad \mu = \begin{pmatrix} 2 \\ 0.8 \\ -0.7 \end{pmatrix}; \quad \Sigma = \begin{pmatrix} 0.7 & 0.8 & 0.8 \\ 0.8 & 1.2 & 0.8 \\ 0.8 & 0.8 & 1.0 \end{pmatrix}.$$

correlation matrix for the time-dependent covariates:

$$\rho = \begin{pmatrix} 1 & 0.873 & 0.956 \\ 0.873 & 1 & 0.730 \\ 0.956 & 0.730 & 1 \end{pmatrix}.$$



- ▶ Data from a Belgian financial institution and add **macro-economic factor**-information from the National Bank of Belgium.
- ▶ Loans ran between 2004 and 2014.
- ▶ 20 000 observations, loan term= 36 months.
- ▶ 7 time-fixed covariates are included.
- ▶ We looked at 20 different models (and the model without TVCs), each time including 3 or 4 of the 6 macro-economic factors.

	Description
z_1	Annual income (per 1000)
z_2	age
z_3	Monthly child allowance (Y/N)
z_4	Number of years at current address
z_5	Total employment years
z_6	Application score
z_7	Mortgage on real estate (Y/N)

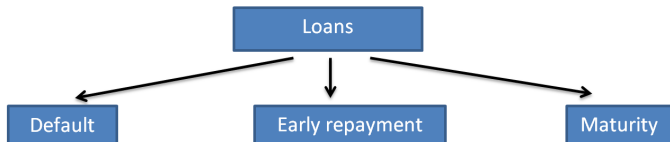
	Description	lag
TVC 1	Market interest rate	6
TVC 2	BEL 20 index	0
TVC 3	Consumer confidence	0
TVC 4	Gross Domestic product	0
TVC 5	Inflation rate	6
TVC 6	Unemployment rate	6

AICcd	(int)	z ₁	z ₂	z ₃	z ₄	z ₅	z ₆	z ₇	IR	BEL 20	cons conf	GDP	inflation	
no TVC	b	2.492 (0.015) ***	0.016 (0.002) ***	0.005 (0.001) ***	-0.076 (0.033) **	-0.023 (0.001) ***	-0.039 (0.002) ***	-1.083 (0.038) ***	-1.319 (0.026) ***					
30127	β		-0.000 (0.005) ns	-0.011 (0.003) ***	-0.074 (0.077) ns	-0.024 (0.004) ***	-0.023 (0.006) ***	-0.615 (0.104) ***	-0.488 (0.092) ***					
best	b	3.55 (0.085) ***	0.046 (0.002) ***	0.058 (0.002) ***	0.192 (0.037) ***	0.009 (0.002) ***	-0.07 (0.002) ***	-1.406 (0.04) ***	-2.153 (0.033) ***					
26685	β		-0.011 (0.005) *	-0.026 (0.003) ***	-0.188 (0.076) *	-0.036 (0.004) ***	-0.016 (0.006) **	-0.519 (0.101) ***	-0.119 (0.077) ns		-0.003 (0.005) ns	-0.011 (0.021) ns	0.001 (0.041) ns	
2nd best	b	2.861 (0.079) ***	0.035 (0.002) ***	0.034 (0.001) ***	0.067 (0.036) .	-0.003 (0.002) *	-0.059 (0.002) ***	-1.23 (0.039) ***	-1.829 (0.032) ***					
27744	β		-0.008 (0.005) .	-0.021 (0.003) ***	-0.133 (0.077) .	-0.031 (0.004) ***	-0.016 (0.005) **	-0.543 (0.1) ***	-0.188 (0.075) *	0.071 (0.036) .	-0.14 (0.057) *		0.009 (0.022) ns	-0.033 (0.041) ns

- ▶ Lower AICcd for all models with macro-economic factors.
- ▶ Interest rate and BEL 20 index only significant TVCs.

Given \mathbf{z} , $\mathbf{x}(t)$, default indicator Y_d , early repayment indicator Y_e

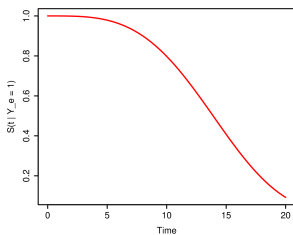
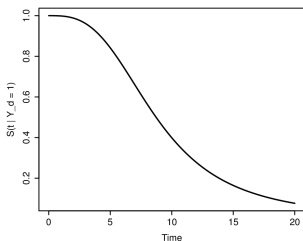
$$S(t | \mathbf{x}) = \pi_d(\mathbf{z})S_d(t | Y_d=1, \mathbf{z}, \mathbf{x}(t)) + \pi_e(\mathbf{z})S_e(t | Y_e=1, \mathbf{z}, \mathbf{x}(t)) + (1 - \pi_d(\mathbf{z}) - \pi_e(\mathbf{z}))$$



prob = $\pi_d(\mathbf{z})$

prob = $\pi_e(\mathbf{z})$

prob = $(1 - \pi_d(\mathbf{z}) - \pi_e(\mathbf{z}))$



<i>d/e</i>		(int)	z_1	z_2	z_3	z_4	z_5	z_6	z_7	IR	BEL 20	GDP	inflation
<i>d</i>	<i>b</i>	-0.033	0.006	-0.014	-0.103	-0.037	-0.047	-1.246	-1.208				
		0.163	0.004	0.003	0.078	0.004	0.005	0.1	0.083				
	β	ns	ns	***	ns	***	***	***	***				
		-0.001	-0.001	-0.049	-0.007	-0.006	-0.198	-0.107	0.039	0.011	0.007	0.029	
		0.005	0.003	0.079	0.004	0.006	0.103	0.086	0.038	0.061	0.022	0.042	
		ns	ns	ns	.	ns	.	ns	ns	ns	ns	ns	
<i>e</i>	<i>b</i>	0.695	-0.007	-0.017	0.008	-0.017	-0.014	-0.741	-0.217				
		0.095	0.002	0.001	0.034	0.002	0.002	0.041	0.034				
	β	***	***	***	ns	***	***	***	***				
		0.002	-0.001	0.037	-0.002	-0.001	-0.176	-0.136	0.01	-0.056	0.023	0.022	
		0.002	0.001	0.036	0.002	0.002	0.042	0.035	0.019	0.025	0.009	0.023	
		ns	ns	ns	ns	ns	***	***	ns	*	*	ns	

- ▶ More complex model: significance of some parameters affected
- ▶ In general, early repayment affected by BEL 20 index and GDP growth

- ▶ AICcd showed that including TVCs can lead to better model fit.
- ▶ Only limited number of TVCs are significant, eg. unemployment. insignificance may be due to selection bias.
- ▶ Information might be lost because monthly averaging: look at weekly default and TVCs.

- ▶ Dirick, L., Claeskens, G. and Baesens, B. (2015). An Akaike information criterion for multiple event mixture cure models. *Eur. J. Oper. Res.*, **241**:449–457.
- ▶ Tong, E. N. C., Mues, C., and Thomas, L. C. (2012). Mixture cure models in credit scoring: if and when borrowers default. *Eur. J. Oper. Res.*, **218**, 132–139.
- ▶ Cavanaugh, J. E. and Shumway, R. H. (1998). An Akaike information criterion for model selection in the presence of incomplete data. *J. Stat. Plan. Inference*, **67**, 45–65.
- ▶ Stepanova, M. and Thomas, L. C. (2002). Survival analysis methods for personal loan data. *Operations Research*, **50 (2)**, 277–289.
- ▶ Watkins, J. G. T., Vasnev, A. L., and Gerlach, R. (2013). Multiple event incidence and duration analysis for credit data incorporating non-stochastic loan maturity. *J. Appl. Econ.*
- ▶ Bellotti, T. and Crook, J. (2009). Credit scoring with macroeconomic variables using survival analysis. *J. Oper. Res. Soc.*, 60(12):1699-1707.
- ▶ Andersen, P. K. (1992). Repeated assessment of risk factors in survival analysis. *Stat. Methods Med Res*, **1(3)**, 297–315.
- ▶ Austin, P. C. (2012). Generating survival times to simulate cox proportional hazards models with time-varying covariates. *Stat. Med.*, **31(29)**:3946–3958.