

**Event History Analysis
Survival Analysis
Duration Analysis
Transition Data Analysis
Hazard Rate Analysis**



Stiftelsen Frischsenteret for samfunnsøkonomisk
forskning
Ragnar Frisch Centre for Economic Research
www.frisch.uio.no

Why Event History Analysis?

- A general analytical framework: The analysis of *discrete events/decisions over time*.
- The method exploits the *timing of events*, not only their occurrence.
- Can be used to assess impacts on events/decisions of:
 - Exogenous covariates
 - Time – Duration – Past events and durations
 - Exogenous and endogenous events
- Particularly suited for efficiently exploiting the wealth of information embedded in merged administrative register data.

Two Interpretations

1. Interpret the dependent variable(s) as a *duration* or (if more than one event can occur) as a set of *latent durations*.
2. Interpret the dependent variable(s) as a sequence of dichotomous (0/1) variables (panel data).

Examples

- Time to recovery (or death) after a medical treatment
- Time to a car accident occurs
- Time to promotion (or dismissal) after first employment
- Time to employment (or labor market exit) after entry into unemployment or after school graduation
 - The change in time to employment caused by participation in an active labor market program (ALMP)
- The duration of poverty
 - The change in the occurrence and duration of poverty caused by activation (or other "treatments")
- The duration of strikes
- Time to retirement
 - The change in time to retirement caused by a change in retirement incentives (e.g., the introduction of AFP)
- The duration of sickness absence
 - The change in the duration of sickness absence caused by a change in certification rules (i.e., July 2004)
 - The impact of absence duration on the relapse propensity
- The length of schooling
- The duration of marriage
 - The change in the duration of marriage caused by child birth

Why Not Standard Regression?

- Regressions with a duration as the dependent variable run into problems due to:
 - Right censoring – not all subjects experience the event(s) in question
 - Time varying right-hand-side variables
- Standard panel data models may run into difficulties due to
 - Many periods (if the information in the data is to be exploited efficiently)
 - Endogenous exits from the panel

Structural or Reduced Form Modeling

- Event history models are typically applied in a *reduced form* setting, with focus on
 - duration dependence,
 - causal effects.
- But structural interpretations can be accommodated, typically with focus on
 - "Deep structural parameters" (i.e., parameters that characterize preferences and/or technology).
- A structural model imposes restrictions on the interpretation of the data, (presumably) justified by theory.
 - The cost: Less reliable results (particularly if someone questions the empirical relevance of the theory).
 - The gain: More general wisdom can be extracted from the data. Results can more easily be extrapolated outside the environment on which they are based.

Basic Concepts in a Single Risk Model

- Let T be a *stochastic* duration, with a continuous distribution on the positive axis and let t be a realization of this variable.
 - The stochastic duration assumption requires that there is always an element of randomness, *even if we knew the data generating process exactly*. Defective risks are ruled out.
- The following functions constitute the building blocks of event history analysis:
 - The distribution function: $F(t) = P(T \leq t)$
 - The survival function: $S(t) = P(T \geq t) = 1 - F(t)$
 - The density: $f(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t)}{\Delta t} = F'(t) = -S'(t)$
 - The hazard rate: $h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$

The Hazard and the Survival Functions

- According to the law of conditional probabilities we have:

$$f(s) = h(s)S(s)$$

$$\Rightarrow h(s) = \frac{f(s)}{S(s)}$$
 - By definition:

$$f(s) = -S'(s)$$

Division by $S(s)$:

$$\frac{f(s)}{S(s)} = \frac{-S'(s)}{S(s)} = h(s)$$
- Integration from 0 to t :
- $$-\int_0^t \frac{S'(s)}{S(s)} ds = \int_0^t h(s) ds$$
- $$\Rightarrow -(\ln S(t) - \underbrace{\ln S(0)}_{=0}) = \int_0^t h(s) ds$$
- $$\Rightarrow \ln S(t) = -\int_0^t h(s) ds$$
- $$\Rightarrow S(t) = \exp \left(- \underbrace{\int_0^t h(s) ds}_{\text{The integrated hazard}} \right)$$

Expected Duration

- Expected duration

$$E(T) = \int_0^{\infty} tf(t)dt$$

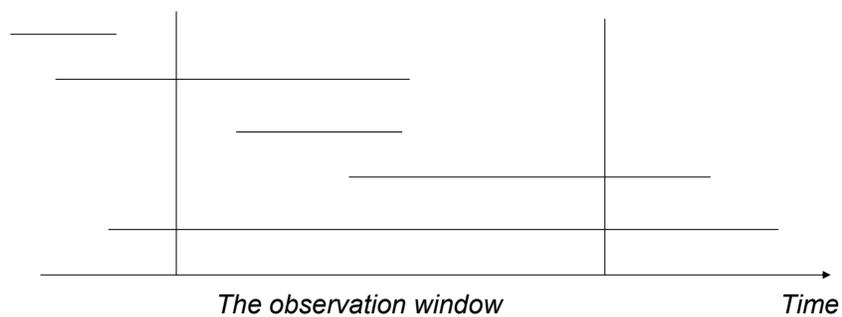
$$= \int_0^{\infty} S(t)dt$$

If the hazard rate is constant (θ):

$$E(T) = \frac{1}{\theta}$$

Data Structures

- Stock or flow sampling?



The Difference Between Stock and Flow Populations

- Example: Mean unemployment duration (based on interviews with stock and “flow out” samples, Norway 1999).
 - Mean **ongoing duration** in stock: 357 days. If the spells on average are sampled halfway, these spells will on average have lasted around 700 days when completed.
 - Mean **completed duration** in “flow out”: 214 days.
- Stock-samples are **length biased**

Censoring and Truncation

- **Censoring**: The researcher observes only a part of a spell, and knows that the duration is *at least as long* as the observed part.
 - Left censoring
 - Right censoring
 - Independent (e.g., end of observation window)
 - Dependent (competing risk – non-random attrition)
- **Truncation**: Some spells (typically very short) are not observed at all (left-truncation).

Nonparametric Descriptions of the Survival Function

- Hazard rates computed on the basis of observed event frequencies in relation to the **risk set**, i.e., the number of subjects at risk at any point in time.
- Life Tables
 - Based on ex ante definition of discrete *time intervals*
- Product Limit Estimator (Kaplan-Meier)
 - Time intervals decided on the basis of actual events

Parametric Specification of Hazard Rates

- Proportional hazard (PH): $h(t | x) = b(t) \underbrace{\lambda(x)}_{\exp(x' \beta)}$
- Accelerated failure time (AFT): $h(t | x) = b(t \lambda(x)) \lambda(x)$

Interpretation of Parameters

- Independent of how time is measured (days, weeks, months), since the hazard is defined in continuous time.
- When the model is proportional, the coefficients attached to explanatory variables can be interpreted as
 - If the variable is a continuous scalar: The elasticity with respect to $\exp(\text{variable})$. So if the variable is measured in logs, we get the elasticity directly.
 - If the variable is a dummy: The impact of the dummy is to shift the hazard by $100(\exp(\text{coefficient})-1)$ percent (relative to a reference value). If the coefficient is “small”, the percentage shift can be approximated by $100*\text{coefficient}$.

The Baseline Hazard in Proportional Models

- Parametric assumptions regarding the distribution:
 - Exponential distribution: $b(t) = b$
 - Weibull distribution: $b(t) = \alpha t^{\alpha-1}$
- Piecewise constant hazard: $b(t) = b_\tau$ for $\tau - 1 < t \leq \tau$

Estimation – Maximum likelihood

- A subject's contribution to the likelihood is either

$$l(t | x) = f(t | x) = h(t | x)S(t | x)$$

for completed durations, or

$$l(t | x) = S(t | x)$$

for right-censored durations

- So, with outcome indicator y , the likelihood function becomes:

$$L = \prod_{i \in N} h(t_i | x_i)^{y_i} S(t_i | x_i)$$

Time-varying Covariates and Episode-splitting

- If explanatory variables are time-varying, each spell is split into sub-spells, such that all variables are constant within each sub-spell.
- With time and/or duration dummy variables – representing piecewise constant hazards – spells must be split according to the assumed frequency of changes in these variables (e.g., weekly or monthly).
- All spell parts are treated as right-censored, except (potentially) the last one.
- The density for a given spell is the product of the survivor functions for all survived sub-spells times the hazard at the time of the transition (unless the spell is censored).

Cox Partial Likelihood

- If the model is of the proportional hazards type, the impacts of x can be estimated while leaving the baseline hazard unspecified.
- Only the *order* of the transitions is exploited, not their exact *timing*.
- Consider an event occurring at duration t , with N_t subjects in the risk set. The probability that this event belongs to individual i is:

$$P(T_i = t | i \in N_t) = \frac{b(t) \exp(x_i' \beta)}{\sum_{j=N_t} b(t) \exp(x_j' \beta)} = \frac{\exp(x_i' \beta)}{\sum_{j=N_t} \exp(x_j' \beta)}$$

Unobserved Heterogeneity

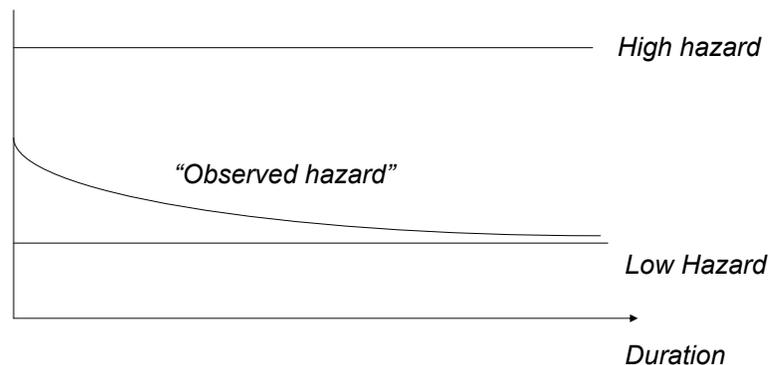
- Mixed proportional hazard (MPH)
 $h(t | x, v) = b(t) \lambda(x) v$
- Since v is unobserved, we cannot use it as an input to the data likelihood function.
- What if we disregard v ?

Unobserved Heterogeneity and Estimation Bias

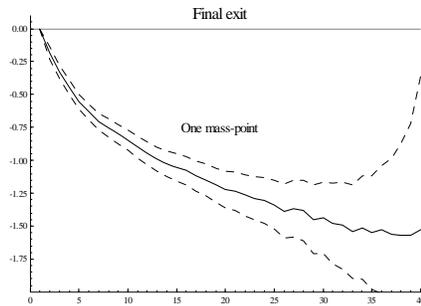
- If v is correlated to x , estimates of β will be biased just like in standard regression analysis (it will measure the causal impact of x as well as the spurious association).
- But even if v is “white noise” (completely random), its existence will bias estimates of both the degree of duration dependence $b(t)$ and the estimated impacts of exogenous covariates β .
- In practice, this implies that is rarely defensible to ignore unobserved heterogeneity.

Sorting and Duration Bias

The nature of sorting by spell duration



Bias Caused by “White Noise” – Evidence from Monte Carlo Trials



Estimated impacts of exogenous variables without control for unobserved heterogeneity Monte Carlo evidence based on 100 trials with 50,000 observations in each trial

	True value	Estimated value	Estimated standard error
Parameter 1	-1	-0.737	0.010
Parameter 2	1	0.968	0.018

How to Take Unobserved Heterogeneity into Account?

- We can write the likelihood contribution (the density) from a completed spell as

$$l(t | x) = E_v[l(t | x, v)] = \int h(t | x, v) S(t | x, v) dv$$

- With outcome indicator y , the likelihood function then becomes:

$$L = \prod_{i \in N} \int h(t_i | x_i, v)^{y_i} S(t_i | x_i, v) dv$$

- But this expression has a tractable (closed form) solution only for some special cases.

The Weibull-Gamma Mixture

- Assume that unobserved heterogeneity is distributed according to a Gamma distribution and that individual spell durations are either exponentially or Weibull distributed.
- Then a closed form solution exists for the marginal distribution of spell durations – the expected survival and density functions can be calculated directly.
- These particular assumptions have often been used simply because they are analytically convenient.
 - They are also incorporated in many software packages.
- **But the results may be completely driven by the assumption!**

Identification

- Unless the impacts of unobserved heterogeneity can be *identified*, it is not meaningful to try to model it.
- The proportionality assumption – combined with at least one exogenous covariate – is sufficient for identification, provided that no risks are *defective*; i.e., all hazards are strictly between 0 and infinity.
- **But can we rely on the proportionality assumption for identification?**
- **After all, the proportionality assumption is typically made for convenience, not because we have prior evidence that it is valid.**

Sources for *Nonparametric* Identification

- Repeated spells
 - Requires that the spells are independent except through the common unobserved variable: No “lagged” duration dependence.
- Exogenous time varying covariates
 - Cyclical or seasonal fluctuations in hazard rates
 - Exogenous events
 - Time variation in the accessibility of treatments

Nonparametric Modeling of Unobserved Heterogeneity – A Latent Class Approach

- Assume that there are Q different “types” – or **latent classes** – v_1, \dots, v_Q , that occur with the probabilities p_1, \dots, p_Q .
- If a subjects contribution to the likelihood conditional on v is $l(t|x, v)$, the expected unconditional contribution is

$$l(t | x) = \sum_{q=1}^Q p_q l(t | x, v_q)$$

- *Nonparametric maximum likelihood (NPMLE): Maximize the likelihood function with respect to all the model parameters plus v_q , p_q , and (NB!) Q .*

How to Select the Number of Classes? Information Criteria for Maximum-Likelihood Estimation

- Maximize the likelihood?
 - Will not the likelihood always increase as more classes are added?
- Maximize a penalized likelihood?
 - Penalized $\log L = \log L - a(\#\text{parameters})$.

BIC: $a = 0,5 \ln(\#\text{obs})$

HQIC: $a = \ln(\ln(\#\text{obs}))$

AIC: $a = 1$

ML: $a = 0$

Interval Censoring

- In practice, we rarely observed the exact timing of an event. Instead, we observe in which month/week/day it occurred.
- We then have three options:
 1. "Pretend" that the data are really continuous (i.e., disregard the problem).
 2. Use a discrete panel data model instead.
 3. Formulate the event history model in terms of interval censored data.

Grouped Hazard Rates

- The probability that an event occurs between $t-1$ and t , given that it did not occur before $t-1$:

$$\begin{aligned}
 p(t-1 < T < t | T > t-1) &= \frac{S(t-1) - S(t)}{S(t-1)} \\
 &= 1 - \frac{S(t)}{S(t-1)} \\
 &= 1 - \exp\left(-\int_{t-1}^t h(s) ds\right) \approx h(t) \text{ for } \int_{t-1}^t h(s) ds \text{ "small"}
 \end{aligned}$$

A complementary log-log function

- Let $h(t, x_t) = b(t) \exp(x_t' \beta)$
- Let x_t be constant within time intervals of unit length.
We can then write

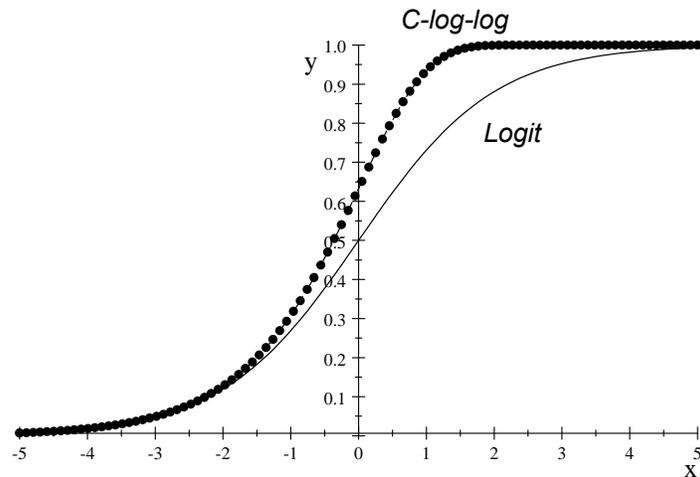
$$\int_{t-1}^t h(s) ds = \exp(x_t' \beta) \int_{t-1}^t b(s) ds = \exp(x_t' \beta + \lambda_t),$$

$$\text{where } \lambda_t = \log \int_{t-1}^t b(s) ds$$

- The probability that an event occurs between $t-1$ and t (given that it did not occur before $t-1$) is then given by the "complementary log-log function" (c-log-log)

$$p(t-1 < T < t | T > t-1) = 1 - \exp\left(-\exp(x_t' \beta + \lambda_t)\right)$$

C-log-log and Logit compared



Frisch Centre 

Interval Censoring and Left Truncation

- Spells that both start and stop between two “observation posts” are not recorded at all.
- Thus, the analysis population is selected, with too few short spells.
- More seriously: If unobserved heterogeneity is present, a stochastic dependence arises between observed and unobserved covariates at the moment of entry into the dataset.

Frisch Centre 

The Unobserved Heterogeneity Distribution Conditional on Being Observed

- Let $f(v)$ be the unconditional distribution (density) of unobserved heterogeneity, and let $S_0(v)$ be the probability of “surviving” the entry interval.
- The conditional distribution of unobserved heterogeneity at entry into the data can then be derived by **Bayes’ Rule**:

$$P(A | B) = \frac{P(B | A)}{P(B)} P(A), \text{ in our case:}$$

$$f(v | \text{observed}) = \frac{S_0(v)}{E[S_0(v)]} f(v)$$



Competing Risks

- A competing risks model is a model with more than one possible outcome/event.
- Competing risks is a trivial model extension insofar as there is
 - No interval censoring or all events take place at the boundary of time-intervals.
 - No unobserved heterogeneity or only unobserved heterogeneity that can be assumed independent across risks
- The different hazard rates (pseudo survival functions) can then be estimated one by one, with transition to competing risks treated as independently right-censored.



Independent Risks

With two competing risks, the hazard rates are defined as

$$h_k(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t, K = k | T \geq t)}{\Delta t}, \quad k = 1, 2$$

The survival function is then

$$S(t) = \exp\left(-\int_0^t (h_1(s) + h_2(s)) ds\right)$$

The two "pseudo" survival functions are

$$S_k(t) = \exp\left(-\int_0^t h_k(s) ds\right)$$

Hence $S(t) = S_1(t)S_2(t)$, and the probability of observing a transition to, e.g., state 1 at time t is $h_1(t)S_1(t)S_2(t)$

The parameters affecting state 1 transitions can thus be maximised independent of the parameters affecting state 2 transitions.

Dependent Risks

- Dependencies between the two risks arise if
 - The data are interval censored (grouped)
 - Unobserved heterogeneity is correlated across the competing states.
 - No longer necessarily the case that disregarding unobserved heterogeneity yields a negative bias in duration dependence.
- The competing hazards must then be modeled simultaneously.
- With grouped hazards, we need an additional assumption regarding the evolution of the competing hazards within the censored time intervals, e.g., that the hazards are constant within these (short) time periods.

Period-specific Event Probabilities (grouped hazards) with K Competing Risks

Let $\int_{t-1}^t h_k(s) ds = h_{kt}$ (constant state k hazard within time-interval)

$$\begin{aligned}
 & p(t-1 < T < t, K = k | T > t-1) \\
 &= \int_{t-1}^t \left(h_{kt} \exp \left(- \sum_k \int_{t-1}^u h_{kt} ds \right) \right) du \\
 &= \int_{t-1}^t \left(h_{kt} \exp \left(- \sum_k (u - (t-1)) h_{kt} \right) \right) du \\
 &= \left(1 - \exp \left(- \sum_k h_{kt} \right) \right) \frac{h_{kt}}{\sum_k h_{kt}},
 \end{aligned}$$

A Useful Decomposition

We then have that

$$\begin{aligned}
 & p(t-1 < T < t | T > t-1, K = k) \\
 &= \underbrace{\left(1 - \exp \left(- \sum_k h_{kt} \right) \right)}_{\text{Probability that one event occurs}} \times \underbrace{\frac{h_{kt}}{\sum_k h_{kt}}}_{\text{Conditional probability that the event is of type k}}
 \end{aligned}$$

The Multivariate Mixed Proportional Hazard Rate Model (MMPH)

$$h_k(t | x_t, v) = b_k(t) \lambda_k(x_t) v_k, \quad k = 1, 2, \dots, K$$

where

(v_1, \dots, v_k) is subject to a joint distribution

Nonparametric Modeling of Unobserved Heterogeneity in Competing Risks Models

- Example: Two possible events, with latent variables (v_1, v_2)
- Standard practice: Expand the model with a new potential location of each variable and estimate the probabilities of all possible combinations. As a result the number of “types” (probabilities) is expanded according to the pattern $1, 2^2=4, 3^2=9, 4^2=16 \dots$. This becomes computationally intractable, particularly if there are many competing risks. With 5 competing risks, the model expands according to $1, 32, 243, 1024, \dots$
- A better approach: Expand the model with only one new probability at each step, i.e., with a new *vector*.
- Or: Reduce “dimensionality” by means of *factor loading*.

The Impact of Endogenous Events – The Timing of Events (ToE) Approach

- Example: The treatment effect of Active Labor Market Program (ALMP) participation on the employment hazard.
- Sample unemployment spells at inflow. The spells start with a competing risk situation: Transitions to
 - Work
 - Treatment (ALMP)
- Transitions to ALMP causes an immediate shift in the employment hazard (on-program effect)
- If the program ends without a job-transition, a new shift occurs (post-program effect).

Identification of Treatment Effect Requires "No-anticipation"

- The agents are assumed not to react to private information regarding the timing of a forthcoming treatment.
- But they are allowed to respond to knowledge regarding the treatment probability as reflected in the statistical model.

Monte Carlo Analysis

(Gaure, Røed, Zhang, 2007)

- Artificial data designed to resemble "typical unemployment spells in administrative register data" (50.000 spells)
- True model (DGP):
 - No duration dependence
 - No treatment effects
 - Exogenous *time* variation in both hazards. Also exogenous variation du to an observed characteristic.
 - Extensive unobserved heterogeneity with positive sorting into treatment (the two unobservables are positively correlated).



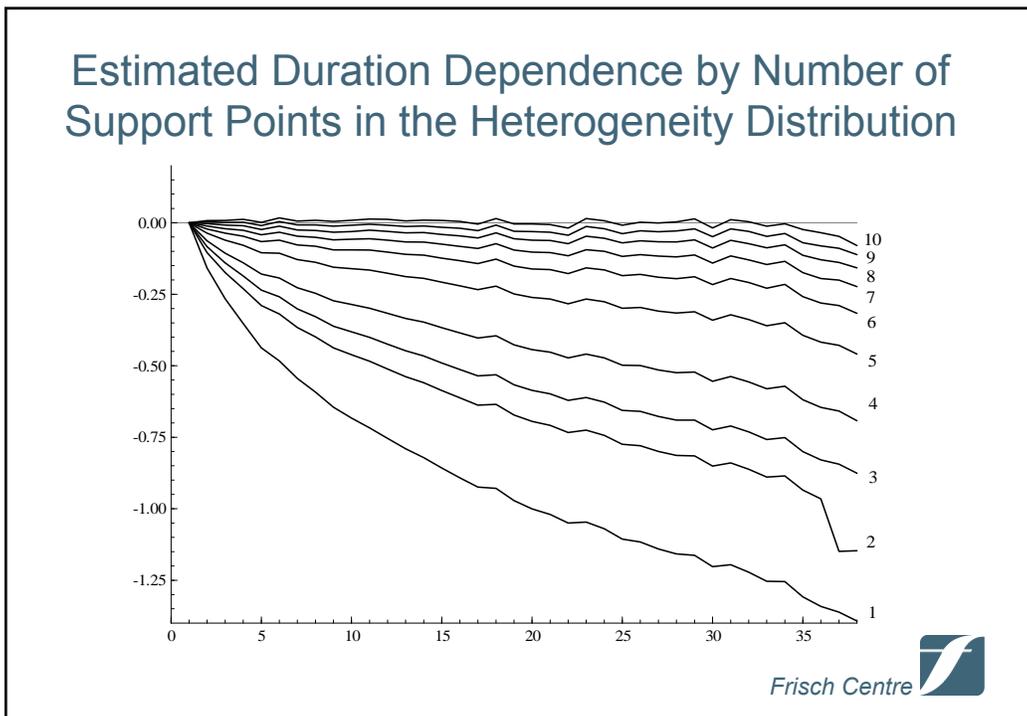
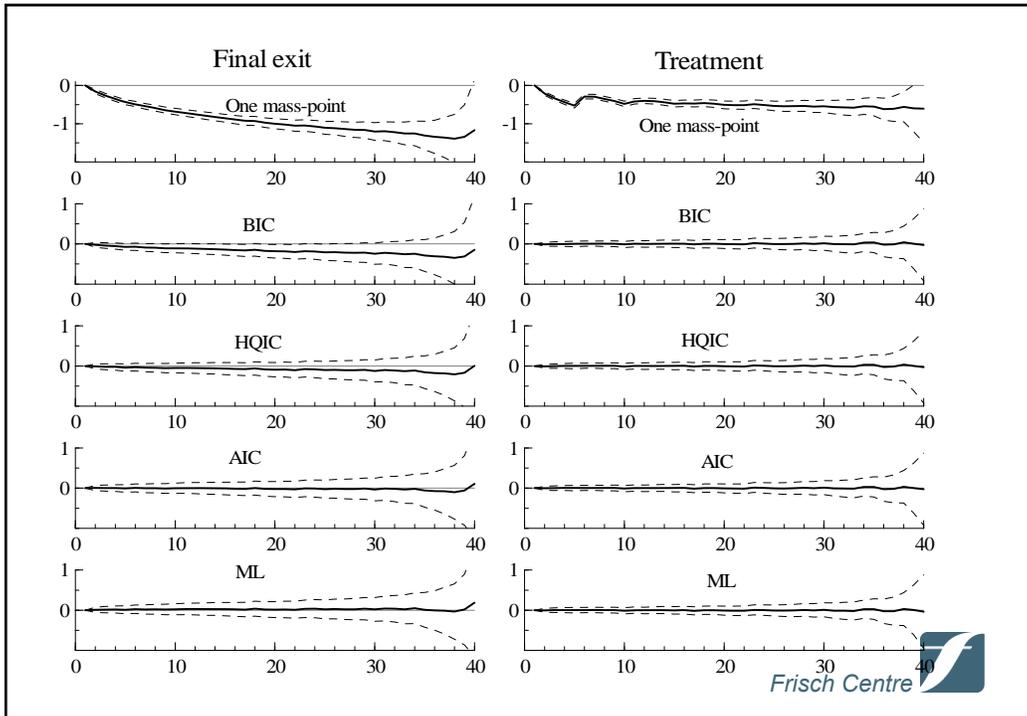
Some Key Results

Table 4
Estimated Effects of Exogenous Covariate and Endogenous Treatment
Results from 100 trials based on the baseline DGP

	True value	Without control for unobserved heterogeneity			BIC			HQIC			AIC			ML		
		Mean Est.	Mean S.E.	Reject at 5%	Mean Est.	Mean S.E.	Reject at 5%	Mean Est.	Mean S.E.	Reject at 5%	Mean Est.	Mean S.E.	Reject at 5%	Mean Est.	Mean S.E.	Reject at 5%
β_z	-1	-0.828	0.014	100	-0.937	0.021	60	-0.968	0.024	31	-0.992	0.026	11	-1.00	0.028	8
β_p	1	0.926	0.011	100	0.989	0.024	8	0.993	0.019	5	0.998	0.026	4	0.998	0.019	3
α_1	0	0.400	0.019	100	-0.020	0.035	19	-0.018	0.037	6	-0.008	0.038	6	-0.003	0.039	4
α_2	0	0.306	0.025	100	-0.022	0.040	21	-0.010	0.039	11	-0.008	0.044	9	-0.003	0.046	6
$\lambda_{z d}$				98			41			20			9			7
$\lambda_{p d}$				96			6			4			5			5

Parameter estimates approximately normally distributed. Standard errors computed conditional on the number of support points can be used for statistical inference.



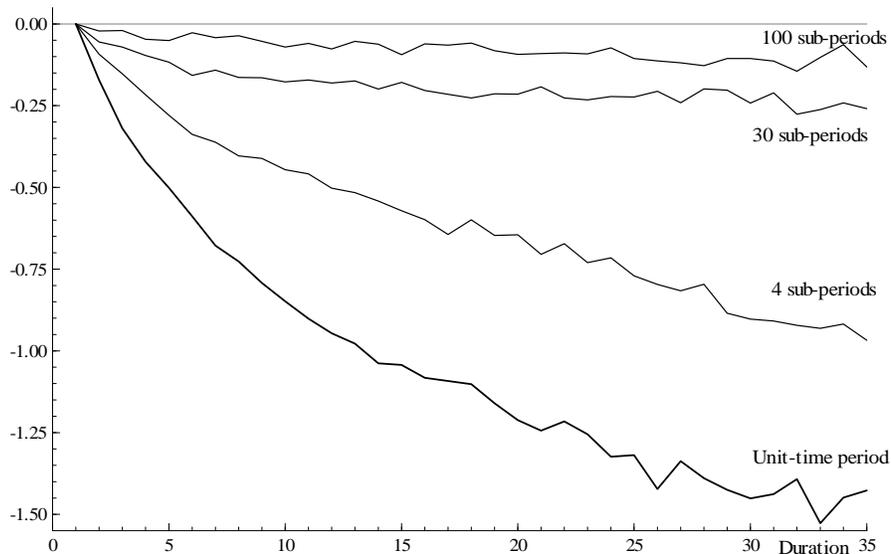


What if We Disregard Interval Censoring?

Table 7
Estimated Effects of Exogenous Covariate and Endogenous Treatment
with Continuous-Time Model Applied on Discrete-Time Data
Results from 10 trials of each model.

	I Integer periods only (40 periods)	II 4 sub-periods between integers (160 periods)	III 30 sub-periods between integers (1,200 periods)	IV 100 sub-periods between integers (4,000 periods)					
Average sub- period event probability (from origin state)	0.186	0.050	0.007	0.002					
Average number of support points found according to ML criterion	4.1	5.8	9.1	9.6					
Estimated effects (and mean standard error)									
	True value	Mean est.	Mean S.E.	Mean est.	Mean S.E.	Mean est.	Mean S.E.	Mean est.	Mean S.E.
β_s	-1	-0.726	0.010	-0.829	0.017	-0.948	0.023	-0.978	0.024
β_p	1	0.951	0.011	1.032	0.185	1.015	0.019	0.997	0.019
α_1	0	0.447	0.026	0.251	0.039	0.041	0.052	0.009	0.042
α_2	0	0.300	0.031	0.206	0.040	0.040	0.044	0.006	0.045

Frisch Centre 



Frisch Centre 

Interpreting the Results - Simulations

- It may be difficult to interpret the output from event history models with many events and many explanatory variables. The coefficient estimates may be of limited interest.
- Model simulations may be used to evaluate the “total impacts” of particular variables or treatments.

Simulation Strategy

- Given vectors of
 - Observed explanatory variables,
 - A distribution of unobserved covariates,
 - A vector of estimated parameters.
- Assign correct observed variables and make draws from the heterogeneity distribution for each subject at entry into the dataset.
- Simulate everything from there (except the path of exogenous covariates)
 - Compare computed transition probabilities over discrete time intervals with draws from a uniform distribution defined on $[0, 1]$,
 - Or invert the (pseudo) survivor function and make a similar draw to predict duration until the event in question occurs.
- Obtain confidence intervals on effects of interest by means of parametric bootstrap.

Parametric Bootstrap

- Draw alternative parameter vectors from the (presumed) multivariate distribution of parameter estimates.
- The draws are made by means of the Cholesky decomposition. Let L be a lower triangular matrix, such that the estimated covariance matrix is $V = LL'$
- Let z_s be a vector of draws from the standard normal distribution collected for trial s . Let \hat{b} be the vector of point-estimates. The parameters drawn for trial s are then given as $b_s = \hat{b} + Lz_s$

Example:

The Estimated Impacts of ALMP Participation

(Gaure, Røed, and Westlie, 2008)

- Topic: The impacts of ALMP on the outcome and duration of unemployed job search and on the quality of a subsequent job.
- Timing of events approach
- The paper identifies:
 - “Adverse” lock-in effect on employment propensity
 - “favorable” post-program effect on employment propensity,
 - Small favorable earnings effect,
 - Small negative effect on job stability.
- **Total impact?**
- Model simulations:
 - Compare outcomes for participants and non-participants *with and without and without treatment*.

	Non-participants	Participants without treatment	Participants with treatment	Effect of treatment
Mean duration of unemployment	5.19	13.95	15.18	1.23 [1.04, 1.41]
Outcomes				
Employment	55.69	47.25	49.32	2.07 [1.46, 2.79]
Education	25.72	25.10	23.52	-1.58 [-2.15, -0.93]
Other Benefits	18.16	25.03	24.98	-0.05 [-0.70, 0.60]
If employment				
Monthly earnings	27,967	25,265	25,908	642 [288, 1,043]
Employment ended within a year	29.63	35.12	36.54	1.42 [0.40, 2.58]