

Curs bàsic d'Anàlisi de dades amb Stata

1

Contingut

- Sessió 4
 - **Regressió lineal**
 - Correlació
 - Regressió Lineal
 - Diagnòstics de regressió
 - Revisió de comandaments d'estimació i post estimació
 - Uso de variables categòriques
 - **Regressió logística**
 - Introducció a la regressió logística
 - Estimació del model
 - Interpretació dels resultats (OR)
 - Confusió e interacció
 - Diagnòstic del model
 - Estratègies de construcció de models de regressió
 - **Anàlisi de Supervivència**
 - Preparació de dades de supervivència: stset
 - Anàlisis descriptiu de dades de supervivència
 - Estimador de Kaplan-Meier
 - Estimació de la funció de Risc
 - Gràfics de supervivència
 - Ajust del model de Cox
 - Diagnòstic del model de Cox
 - **Exercici pràctic**

2

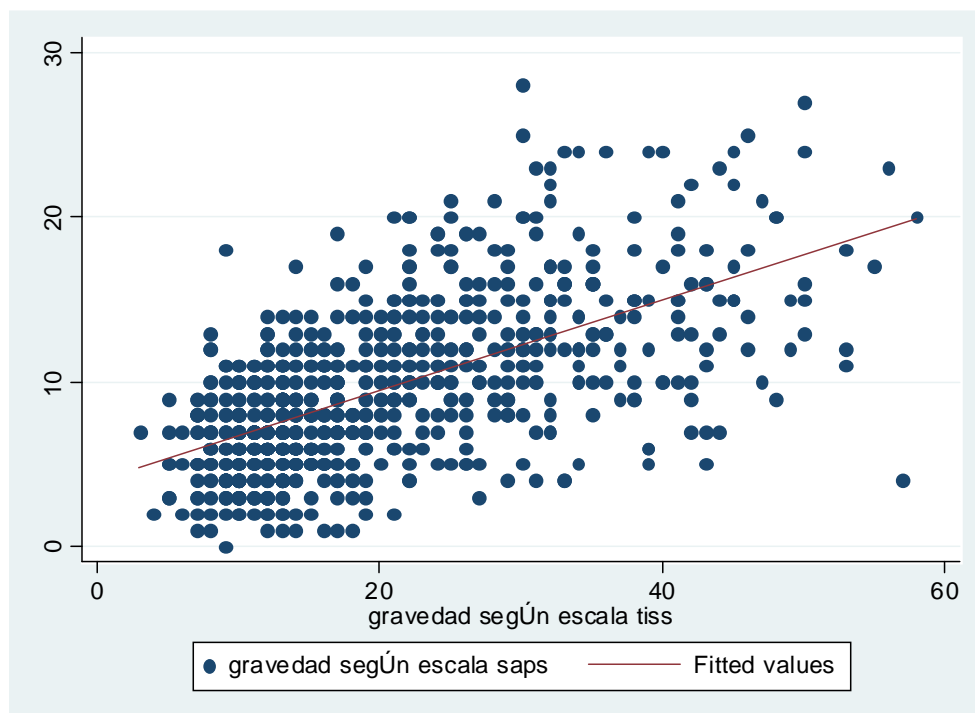
Análisis multivariante

Respuesta	Exposición	Modelos de regresión
Continua	Continua	Regresión lineal Regresión no paramétrica Regresión polinómica o fraccional
Continua	Categórica	Anova , Ancova, Regresión lineal
Dicotómica	Continua	Logística ordinaria o condicional
Dicotómica	Categórica	Regresión logística/ Modelos log-lineales
Recuento/personas-año	Categórica/Continua	Regresión Poisson
Tiempo a evento	Categórica/Continua	Regresión de Cox o modelos paramétricos de supervivencia

3

Regresión lineal con Stata

```
twoway scatter depvar indepvar ||lfit depvar indepvar
twoway scatter saps tiss|| lfit saps tiss
```



4

Regresión lineal con Stata

regress *depvar* [*indepvars*] [*if*] [*in*] [*weight*] [, *options*]

regress *saps* *tiss*

Source	SS	df	MS	Number of obs = 828		
Model	7230.58002	1	7230.58002	F(1, 826) = 512.47		
Residual	11654.1688	826	14.1091632	Prob > F = 0.0000		
Total	18884.7488	827	22.8352464	R-squared = 0.3829		
				Adj R-squared = 0.3821		
				Root MSE = 3.7562		
saps	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
tiss	.2755289	.0121711	22.64	0.000	.2516389	.2994189
_cons	3.897909	.2743255	14.21	0.000	3.359452	4.436366

5

Regresión con variables categóricas

xi:regress *depvar* [*i.varcat*], *options*

xi:regress *saps* *i.educacio*

<i>i.educacio</i>	<i>_Ieducacio_1-4</i>	(naturally coded; <i>_Ieducacio_1</i> omitted)				
Source	SS	df	MS	Number of obs = 803		
Model	360.60094	3	120.200313	F(3, 799) = 5.44		
Residual	17662.6618	799	22.1059597	Prob > F = 0.0011		
Total	18023.2628	802	22.4728962	R-squared = 0.0200		
				Adj R-squared = 0.0163		
				Root MSE = 4.7017		
saps	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
<i>_Ieducacio_2</i>	-1.359817	.4164515	-3.27	0.001	-2.177285	-.5423485
<i>_Ieducacio_3</i>	-1.899653	.5061159	-3.75	0.000	-2.893127	-.9061788
<i>_Ieducacio_4</i>	-1.50506	.7431685	-2.03	0.043	-2.963853	-.0462668
_cons	10.44624	.3447452	30.30	0.000	9.769523	11.12295

6

Graficos de residuos

```
predict nomvarres, residuals[rstudent][rstandard]
```

xi:regress saps i.educacio						
i.educacio	_Ieducacio_1-4	(naturally coded; _Ieducacio_1 omitted)				
Source	SS	df	MS	Number of obs = 803		
Model	360.60094	3	120.200313	F(3, 799) = 5.44		
Residual	17662.6618	799	22.1059597	Prob > F = 0.0011		
Total	18023.2628	802	22.4728962	R-squared = 0.0200		
				Adj R-squared = 0.0163		
				Root MSE = 4.7017		
saps	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
_Ieducacio_2	-1.359817	.4164515	-3.27	0.001	-2.177285	-.5423485
_Ieducacio_3	-1.899653	.5061159	-3.75	0.000	-2.893127	-.9061788
_Ieducacio_4	-1.50506	.7431685	-2.03	0.043	-2.963853	-.0462668
_cons	10.44624	.3447452	30.30	0.000	9.769523	11.12295

7

Construcción de modelos

```
estimates store nommodelo [guarda modelo]
estimates replay nommodelo [activa modelo]
estimates stats nommodelo [activa modelo]
lrtest nommodelo [test ajuste modelo]
```

```

xi:regress saps i.educacio if tiss!=.
estimates store mod1
xi:regress saps i.educacio tiss
lrtest mod1

```

Likelihood-ratio test	LR chi2(1) =	378.58
(Assumption: mod1 nested in .)	Prob > chi2 =	0.0000

```
estimates stats
```

Model	Obs	ll(null)	ll(model)	df	AIC	BIC
mod1	796	-2356.947	-2348.861	4	4705.722	4724.441

•

Regressión logística con STATA

COMANDO variables if condición , Opciones del comando

```
[By varlist:]logit var dep [vars indep] if condición ,  
opciones
```

<u>level</u> (#)	Límite de los intervalos de confianza
<u>nocoef</u>	No se muestra tabla de coeficientes
<u>noconstant</u>	Suprime la constante(intercept)
<u>robust</u>	Proporciona estimaciones robustas del intervalo de confianza
<u>cluster</u> (variable)	Variable que identifica a los sujetos y por tanto observaciones repetidas
<u>or</u>	Muestra los OR en lugar de los coeficientes
<u>offset</u> (variable)	Variable que entra con coeficiente 1

```
[By varlist:]logistic var dep [vars indep] if condición ,  
opciones
```

<u>level</u> (#)	Límite de los intervalos de confianza
<u>robust</u>	Proporciona estimaciones robustas del intervalo de confianza
<u>cluster</u> (variable)	Variable que identifica a los sujetos y por tanto observaciones repetidas
<u>offset</u> (variable)	Variable que entra con coeficiente 1
<u>group</u> (#)	Número de cuantiles para agrupar los datos.
<u>all</u>	Proporciona todos los estadísticos

9

Ajuste modelo con STATA

```
. xi:logit mort tiss_20 edad_60 dias_5 i.sitlabor if validos==1
```

```
i.sitlabor      _Isitlabor_1-7      (naturally coded; _Isitlabor_1 omitted)
```

```
note: _Isitlabor_7 dropped due to collinearity
```

```
Iteration 0:    log likelihood = -265.96082
```

```
Iteration 1:    log likelihood = -198.96442
```

```
Iteration 2:    log likelihood = -179.67671
```

```
Iteration 3:    log likelihood = -178.57059
```

```
Iteration 4:    log likelihood = -178.54761
```

```
Iteration 5:    log likelihood = -178.54759
```

```
Logit estimates
```

```
Number of obs   =          695
```

```
LR chi2(6)      =          174.83
```

```
Prob > chi2     =           0.0000
```

```
Pseudo R2      =           0.3287
```

```
Log likelihood = -178.54759
```

mort	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
tiss_20	.1496574	.0147328	10.16	0.000	.1207816	.1785331
edad_60	.0700162	.0163476	4.28	0.000	.0379756	.1020569
dias_5	-.0635884	.0223158	-2.85	0.004	-.1073266	-.0198503
_Isitlabor_2	1.343407	.5485736	2.45	0.014	.2682227	2.418592
_Isitlabor_3	-.321371	.4595417	-0.70	0.484	-1.222056	.5793143
_Isitlabor_4	.8108309	.6516717	1.24	0.213	-.4664221	2.088084
_cons	-2.53193	.3557546	-7.12	0.000	-3.229197	-1.834664

10

Ajuste modelo con STATA. OR

```
. xi:logit mort tiss_20 edad_60 dias_5 i.sitlabor if noentra==0,or  
...
```

mort	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]	
tiss_20	1.161436	.0171112	10.16	0.000	1.128378	1.195462
edad_60	1.072526	.0175332	4.28	0.000	1.038706	1.107447
dias_5	.9383911	.0209409	-2.85	0.004	.8982323	.9803454
_Isitlabor_2	3.832078	2.102177	2.45	0.014	1.307638	11.23003
_Isitlabor_3	.7251542	.3332386	-0.70	0.484	.2946237	1.784814
_Isitlabor_4	2.249777	1.466116	1.24	0.213	.6272425	8.069439

11

Ajuste regresión logística condicional con STATA

*[By varlist:] clogit var dep [vars indep] if condición ,
opciones*

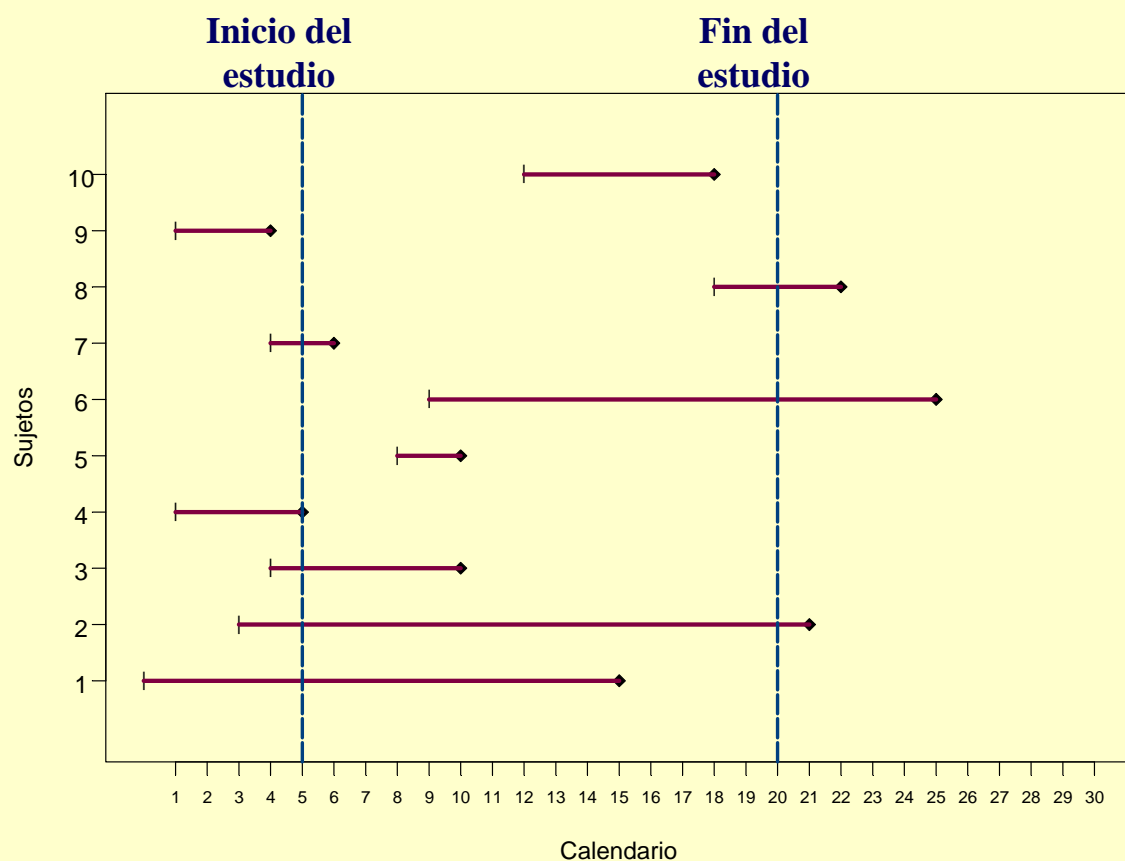
<u>group(variable)</u>	Variable que indica los casos apareados
<u>level(#)</u>	Límite de los intervalos de confianza
<u>or</u>	Muestra los OR en lugar de los coeficientes
<u>offset(variable)</u>	Variable que entra con coeficiente 1

12

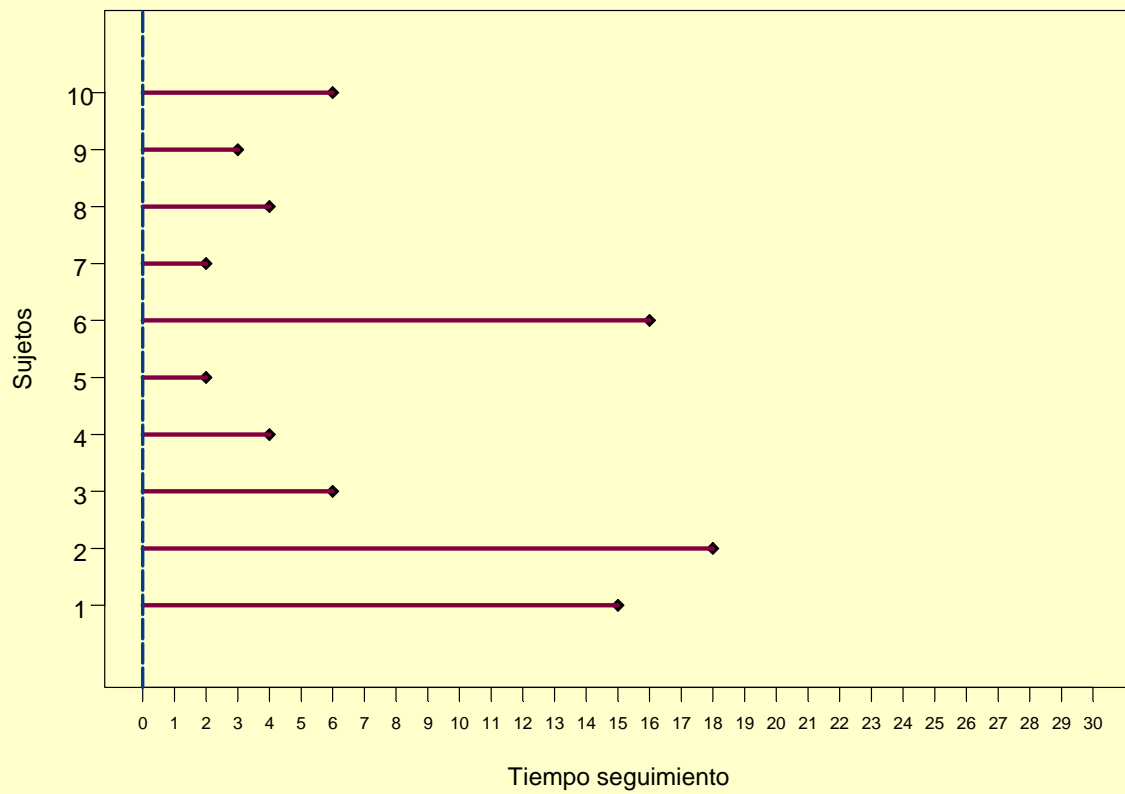
Análisis de Supervivencia con STATA

13

Datos reales tiempo de seguimiento

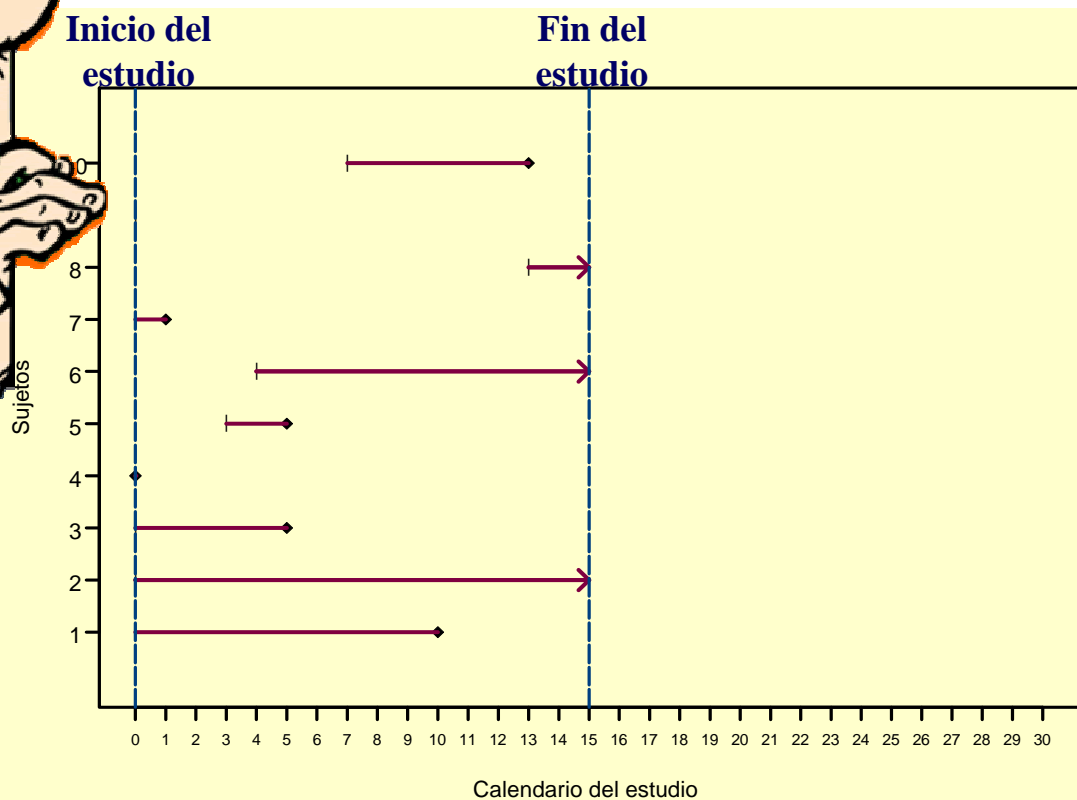


Tiempo de Seguimiento análisis ideal

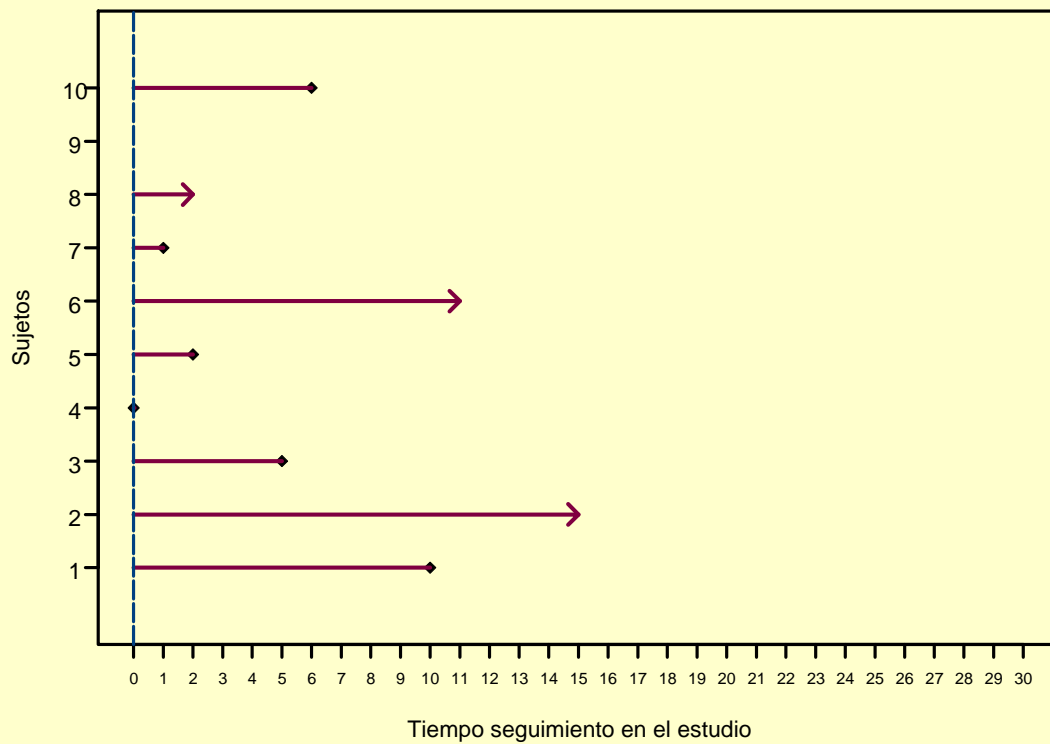


15

Tiempos observados en estudio real

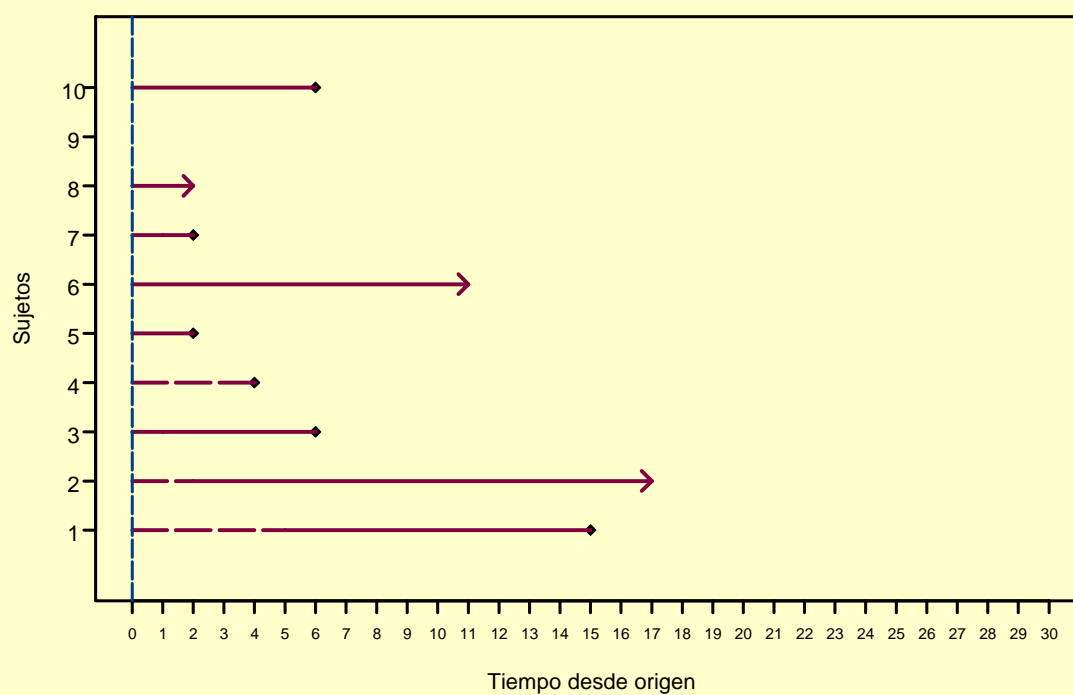


Tiempo de Seguimiento análisis "habitual"



17

Análisis con datos truncados (late entry)



18

Funcion de Supervivencia $S(t)$

$$S(t) = \text{Prob}(\text{Sobrevivir a } t) = P\{T > t\} = 1 - P(\text{fallecer antes de } t)$$

Tasa de peligro $\lambda(t) = h(t)$

Probabilidad de fallecer en un intervalo de tiempo muy pequeño sabiendo que se está vivo al inicio



$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \text{Prob}(\text{fallecer en } (t, t + \Delta t) / \text{vivo en } t) / \Delta t \\ = f(t) / S(t)$$

Percentil t_p

Tiempo en el que el $p\%$ de los sujetos desarrollan el evento

19

Análisis de Supervivencia con STATA

- Todas las instrucciones que tienen que ver con datos de supervivencia van precedidas por el término `st`
- Primer paso: declarar los datos como datos de supervivencia

```
stset  timevar,  failure(event)
time0  (variable) enter(variable)
origin(variable) scale(365,25)
```

20

Análisis de Supervivencia con STATA

Notas:

- Failure() actua como indicador:
0 y missing → “censura”
resto de valores → “muerte”

sujeto	tseg	mort
1	100	1
2	150	0
3	97	.
4	110	1

```
stset tseg, failure(mort)
```

21

Análisis de Supervivencia con STATA

- Si no especificamos failure () : todos los registros acaban en muerte

sujeto	tseg	mort
1	100	1
2	150	1
3	97	1
4	110	1

```
stset tseg
```

22

Análisis de Supervivencia con STATA

- En failure() podemos añadir más de un código

sujeto	tseg	enf
1	100	1
2	150	0
3	97	0
4	110	3

```
stset tseg,failure(enf==1,3)
```

23

Utilizando fechas y diferentes escalas temporales

Escala de tiempo = Tiempo desde el diagnóstico en años:

```
stset dateexit, failure(dead==1) origin(datediag) scale(365.25)
```

datebth	datediag	dateexit	dead	_t0	_t	_d
13feb1906	02jan1986	05feb1986	dead	0	.09308693	1
07mar1906	21jan1986	17feb1986	dead	0	.07392197	1
11mar1906	07jan1986	17jan1986	dead	0	.02737851	1
21apr1906	22jan1986	31jan1986	dead	0	.02464066	1
23apr1906	18feb1986	25jun1986	dead	0	.34770705	1

Escala de tiempo = Edat en años (datos truncados)

```
stset dateexit, failure(dead==1) origin(datebth) enter(datediag) scale(365.25)
```

datebth	datediag	dateexit	dead	_t0	_t	_d
13feb1906	02jan1986	05feb1986	dead	79.88501	79.978097	1
07mar1906	21jan1986	17feb1986	dead	79.876797	79.950719	1
11mar1906	07jan1986	17jan1986	dead	79.827515	79.854894	1
21apr1906	22jan1986	31jan1986	dead	79.756331	79.780972	1
23apr1906	18feb1986	25jun1986	dead	79.824778	80.172485	1

24

Análisis de Supervivencia con STATA

Comandos más importantes

- **stdes** `[if...]` → describe los datos indicados en la instrucción stset
- **stsum** `[if...], by(variables)`
→ presenta estadísticos descriptivos y tasas de incidencia de los datos de supervivencia totales o por los grupos generados por una variable

25

Análisis de Supervivencia con STATA

- **sts list** `[if..], by (variables) failure compare at(instantes de tpo) na`
- muestra estimadores de la supervivencia, su complementario(failure) y la tasa acumulada (na). Se pueden mostrar en unos instantes de tiempo(at) y comparar (compare) en los grupos generados por varias variables (by)

26

Análisis de Supervivencia con STATA

- **sts graph**[if...],by (*variables*)
failure gwood na cna censored
(single/number) lost hazard →
dibuja las curvas de supervivencia de Kaplan-Meier, su complementaria (failure), la tasa acumulada (na). Se pueden dibujar los int. de confianza para la supervivencia (gwood) y para la tasa acumulada (cna). Podemos marcar y enumerar las censuras (censored)
- **sts test variable [if...],**
[logrank/wilcoxon/tware/peto] →
calcula distintos test para comparar la supervivencia entre dos o más grupos

27

Ejemplo cohorte seroconvertidores

```
gen exit_date= datalive
(59 missing values generated)
replace exit_date=dieddate if died==1
(64 real changes made)
stset exit_date,f(mort==1) origin(serodate) scale(365.25)

      failure event:   mort == 1
obs. time interval:   (origin, exit_date]
exit on or before:    failure
t for analysis:       (time-origin)/365.25
      origin:         time serodate

-----
447  total obs.
383  obs. end on or before enter()
-----

64  obs. remaining, representing
64  failures in single record/single failure data
316.5585  total analysis time at risk, at risk from t =          0
              earliest observed entry t =          0
              last observed exit t =    10.80903
```

28

Descriptivo

stdes

```
failure _d: mort == 1
analysis time _t: (exit_date-origin)/365.25
origin: time serodate
```

Category	total	per subject			
		mean	min	median	max
no. of subjects	64				
no. of records	64	1	1	1	1
(first) entry time		0	0	0	0
(final) exit time		4.946227	.2655715	4.724162	10.80903
subjects with gap	0				
time on gap if gap	0				
time at risk	316.55852	4.946227	.2655715	4.724162	10.80903
failures	64	1	1	1	1

29

Descriptivo

stsum, by(expcateg)

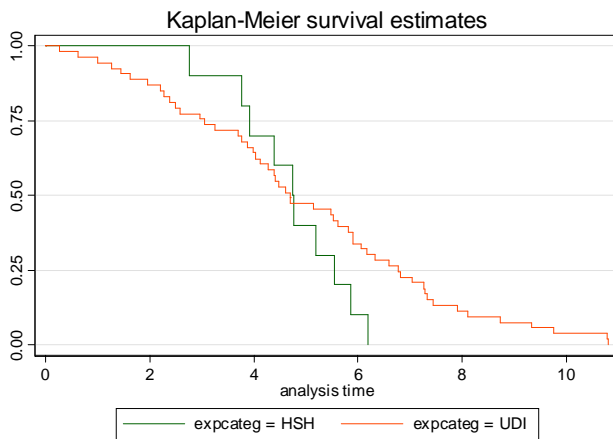
```
failure _d: mort == 1
analysis time _t: (exit_date-origin)/365.25
origin: time serodate
```

expcateg	time at risk	incidence rate	no. of subjects	Survival time		
				25%	50%	75%
HSH	47.15400411	.2120711	10	3.926078	4.73922	5.555099
UDI	264.4982888	.2003794	53	3.058179	4.692676	6.776181
total	311.652293	.2021484	63	3.258042	4.709103	6.329911

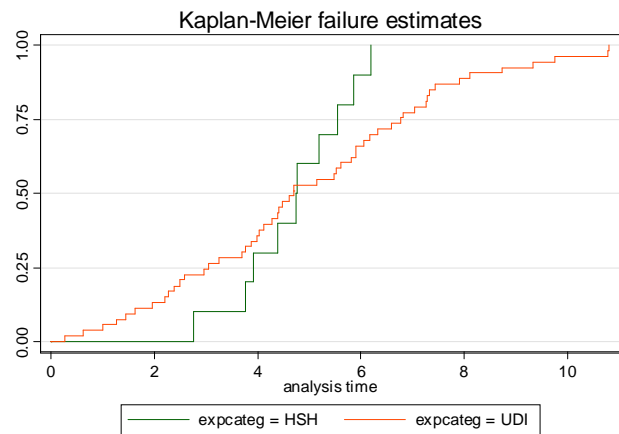
30

Gráficos

sts graph, by(expcateg)



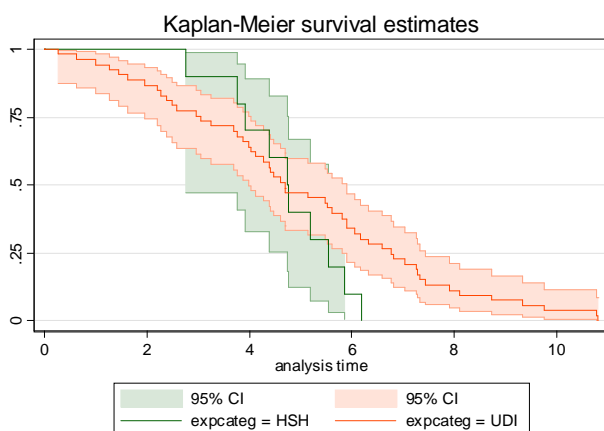
sts graph, by(expcateg) f



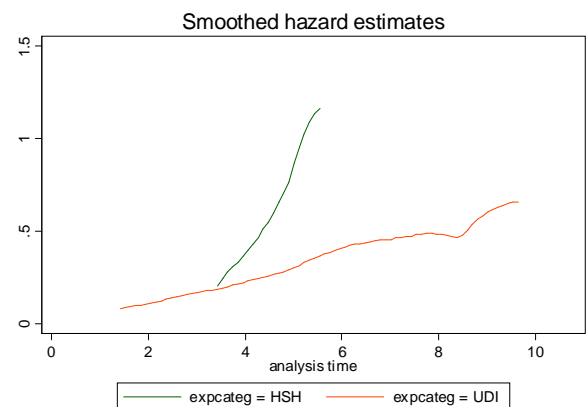
31

Gráficos

sts graph, by(expcateg) ci



sts graph, by(expcateg) h



32

Listado

. sts test expcateg

```
failure _d: mort == 1
analysis time _t: (exit_date-origin)/365.25
origin: time serodate
```

Log-rank test for equality of survivor functions

	Events	Events
expcateg observed	expected	
HSH	10	7.26
UDI	53	55.74
Total	63	63.00

```
chi2(1) = 1.25
Pr>chi2 = 0.2642
```

33

Test Log-rank

sts list, by(expcateg) at(0 1 2 4 6 8 10)

```
failure _d: mort == 1
analysis time _t: (exit_date-origin)/365.25
origin: time serodate
```

	Time	Beg. Total	Fail	Survivor Function	Std. Error	[95% Conf. Int.]	
HSH	0	0	0	1.0000	.	.	.
	1	0	0	1.0000	.	.	.
	2	0	0	1.0000	.	.	.
	4	8	3	0.7000	0.1449	0.3287	0.8919
	6	2	6	0.1000	0.0949	0.0057	0.3581
	8	1	1
	10	1	0
UDI	0	0	0	1.0000	.	.	.
	1	52	2	0.9623	0.0262	0.8574	0.9904
	2	47	5	0.8679	0.0465	0.7428	0.9347
	4	35	12	0.6415	0.0659	0.4973	0.7542
	6	20	16	0.3396	0.0651	0.2168	0.4664
	8	7	12	0.1132	0.0435	0.0460	0.2141
	10	3	4	0.0377	0.0262	0.0070	0.1148

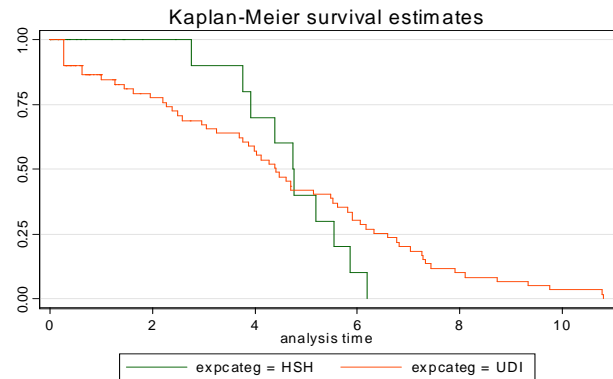
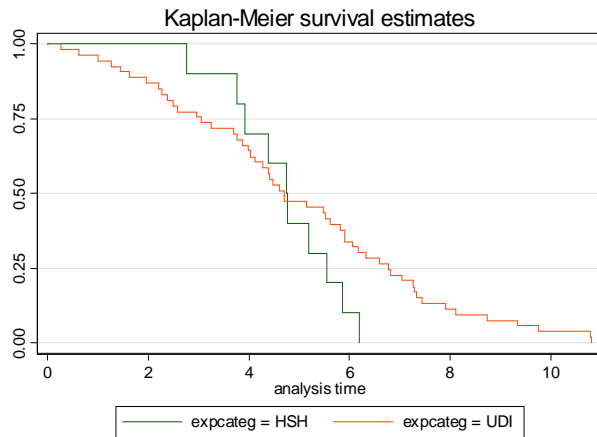
Note: survivor function is calculated over full data and evaluated at indicated times; it is not calculated from aggregates shown at left.

34

Gráficos

```
stset exit_date,  
f(mort==1)  
origin(serodate)
```

```
enter(firstpos)  
scale(365.25)
```



35

Modelo de Cox

$$\log [h(t;x)/h_0(t)] = \beta X$$

Se dispone de datos de la forma (t_i, δ_i, x)

El objetivo es

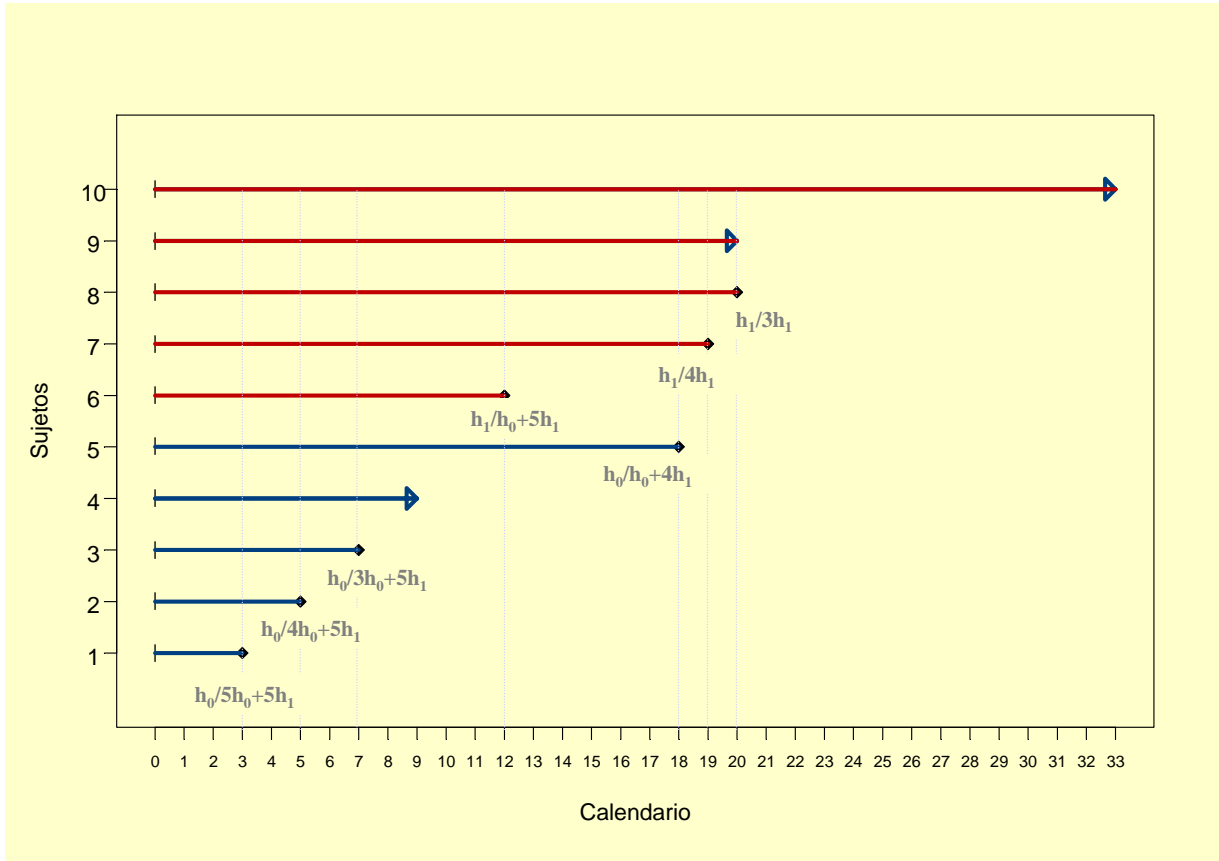
- 1) Estimar β y contrastar la hipótesis $H_0: \beta=0$
- 2) Estimar h_0

Para 1) se utiliza máximoverosimilitud condicional

Para 2) se utilizan métodos no paramétricos

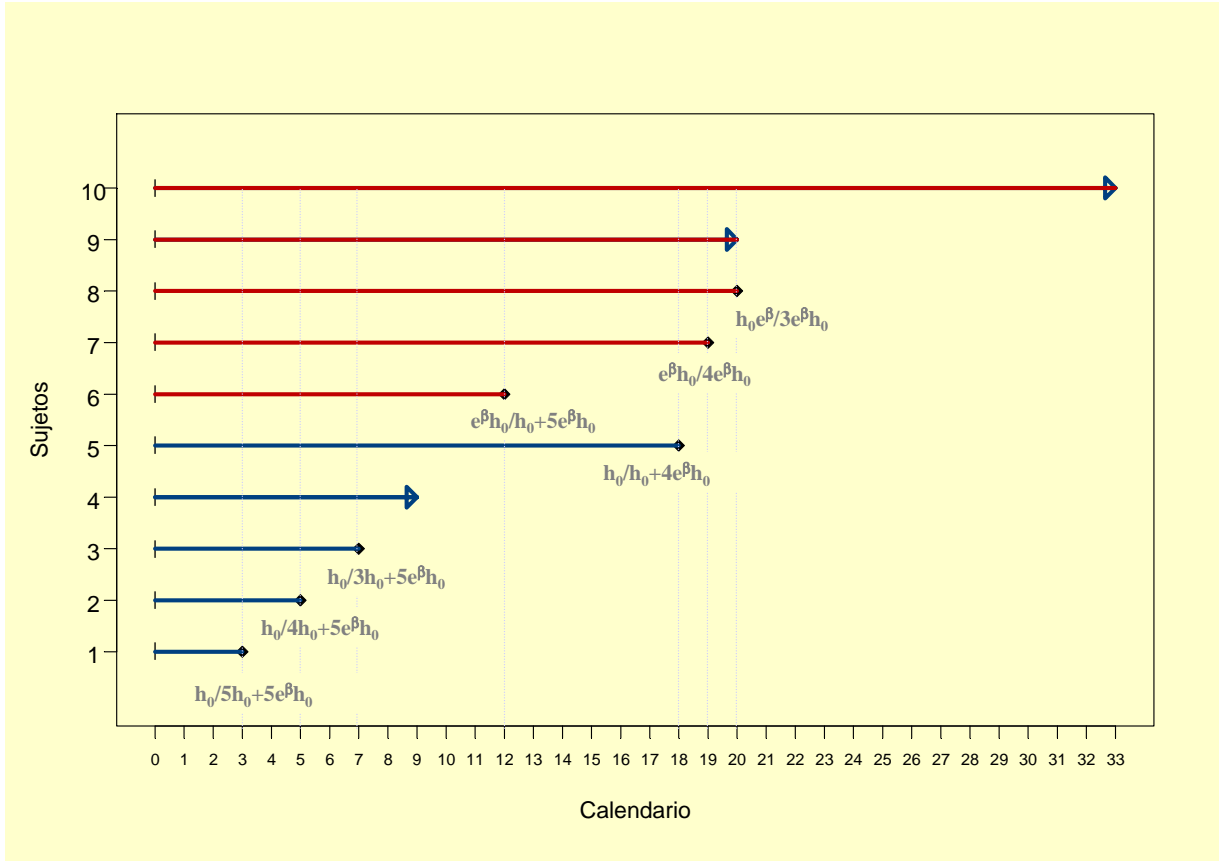
36

Datos Brown



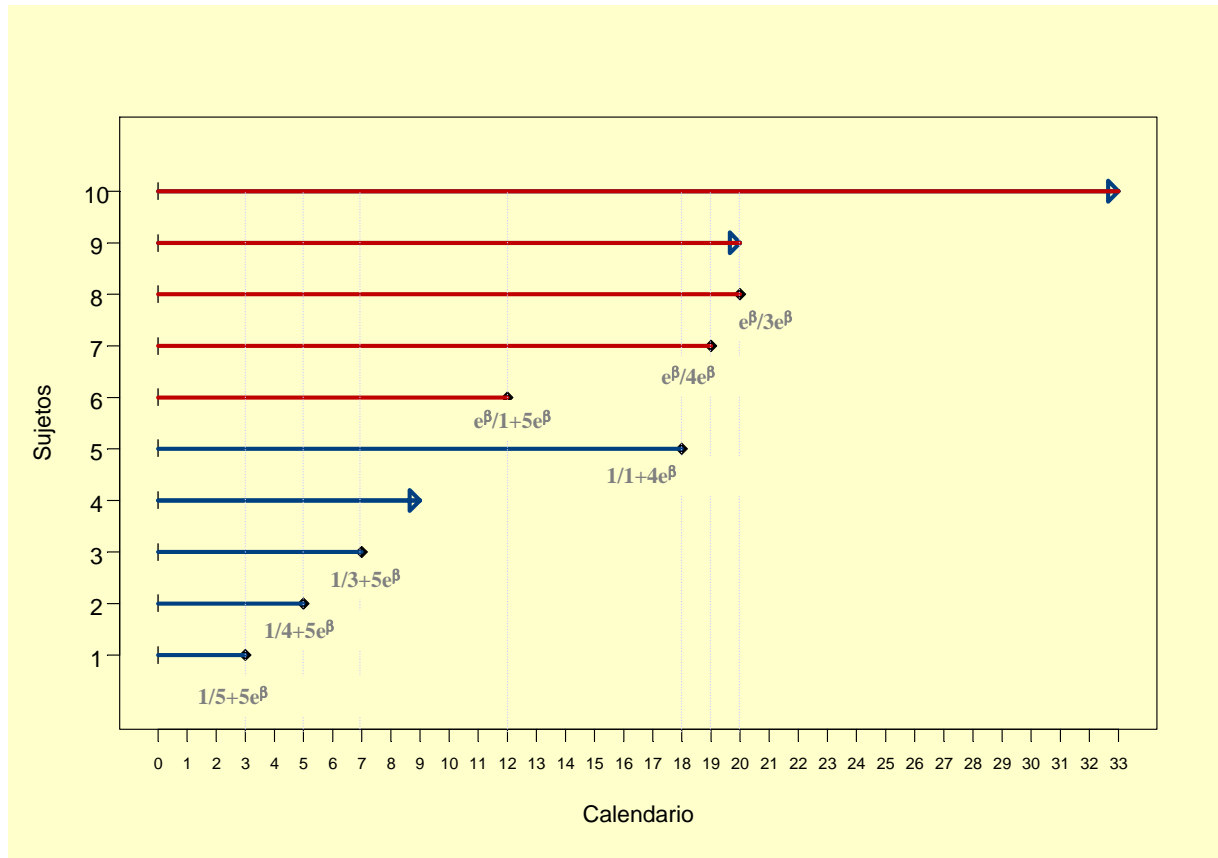
7

Datos Brown



8

Datos Brown



9

Regresión de Cox

Instrucción para la Regresión de Cox

```
xi:stcox [varcontinuas i.varcategóricas][if
exp][in range][,nohr strata(varnames)robust
cluster(varname)noadjust
tvc(varlist)texp(exp)shared(varname)frailty(g
amma)
effects(newvar)mgale(newvar)esr(newvar(s))bas
ehc(newvar)basechazard(newvar)basesurv(newvar
){breslow| efron|exactm|exactp}estimate
noshow
offset(varname)level(#)maximize_options]
```

Regresión de Cox

- **nahr** → muestra los coeficientes en lugar de los hazard ratio
- **robust** → nos da una estimación robusta de la varianza
- **basehazard(newvar)** → añade una nueva variable a los datos que contiene la estimación de la función de peligro (hazard function $H_0(t)$) acumulada
- **basesurv(newvar)** → añade una nueva variable a los datos que contiene la estimación de la función de supervivencia (survival function $S_0(t)$)

41

Regresión de Cox

xi:stcox i.expcateg

i.expcateg _Iexpcateg_1-2 (naturally coded; _Iexpcateg_1 omitted)

failure _d: mort == 1
analysis time _t: (exit_date-origin)/365.25
origin: time serodate

Iteration 0: log likelihood = -201.05811
Iteration 1: log likelihood = -200.49958
Iteration 2: log likelihood = -200.48797
Iteration 3: log likelihood = -200.48797
Refining estimates:
Iteration 0: log likelihood = -200.48797

Cox regression -- Breslow method for ties

No. of subjects =	63	Number of obs =	63
No. of failures =	63		
Time at risk =	311.652293		
Log likelihood =	-200.48797	LR chi2(1) =	1.14
		Prob > chi2 =	0.2856

_t	Haz. Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
_Iexpcateg_2	.6706417	.2416259	-1.11	0.267	.3309868 1.358847

42

Post estimación modelo de Regresión de Cox

43

Ajuste, tests, indicadores

```
estimates store nommodelo      [guarda modelo]  
estimates replay nommodelo    [activa modelo]  
estimates stats nommodelo     [activa modelo]  
lrtest nommodelo              [test ajuste modelo]  
estat concordance             [calcula Harrel's C]
```

```
failure _d:  mort == 1  
      analysis time _t:  (exit_date-origin)/365.25  
              origin:  time serodate
```

Harrell's C concordance statistic

Number of subjects (N)	=	63
Number of comparison pairs (P)	=	1952
Number of orderings as expected (E)	=	274
Number of tied predictions (T)	=	1422

Harrell's C = (E + T/2) / P = .5046
Somers' D = .009221

44

Proporcionalidad a lo largo del tiempo

- *El modelo de Cox asume proporcionalidad a lo largo del tiempo*
- *Por ello al introducir t o $\log(t)$ en el modelo el β correspondiente debería ser 0.*
- *En caso contrario implica que la tasa de peligro depende del tiempo*

45

Proporcionalidad a lo largo del tiempo

```
xi:stcox i.expcateg, tvc(expcateg) texp(ln(_t))
```

```
-----+-----  
      _t | Haz. Ratio   Std. Err.      z    P>|z|     [95% Conf. Interval]  
-----+-----  
main      |  
_Iexpcateg_2 |    56.32939   128.0702     1.77   0.076     .6538049     4853.13  
-----+-----  
tvc        |  
_Iexpcateg_2 |    .0407242    .0605625    -2.15   0.031     .002208     .7511161  
-----+-----
```

Note: variables in tvc equation interacted with ln(_t)

46

Residuos Cox-Snell

- Hay que recordar la relación existente entre la tasa acumulada y la supervivencia $H(t) = -\log(S(t))$
- Supongamos que se sustituyen los tiempos de supervivencia t_i por la tasa acumulada en ese punto y calculemos cual es la supervivencia para la variable aleatoria $H(t_i)$

$$P(H_i(T) > t) = \quad (\text{aplicando la función inversa})$$

$$P(T > H_i^{-1}(t)) = S(H_i^{-1}(t)) \quad (\text{aplicando relación entre } S(t) \text{ y } H(t))$$

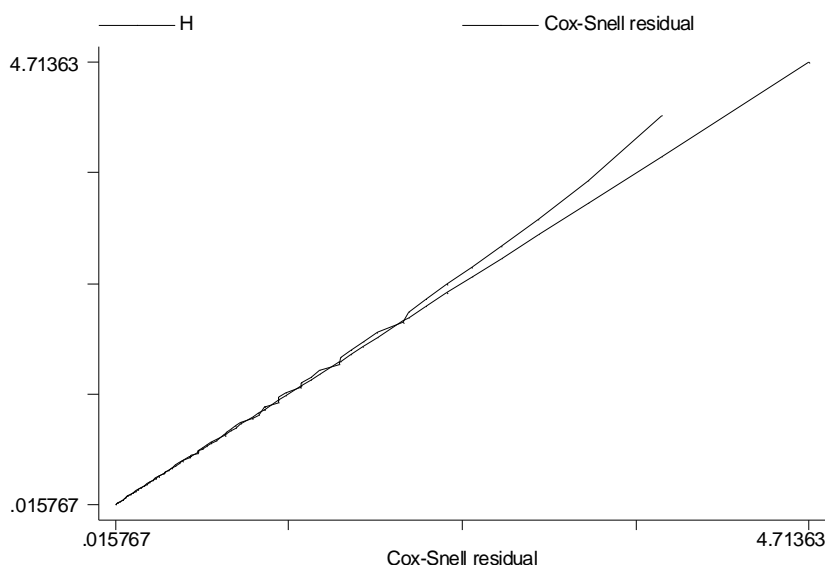
$$= \exp(-H_i(H_i^{-1}(t))) = \exp(-t)$$

- Así $(H_i(t_i), \delta_i)$ son una muestra de datos censurados que siguen una distribución exponencial de media 1.
- Si el modelo de riesgos proporcionales es adecuado $(\exp(\beta x_i) H_0(t_i), \delta_i)$ deben de seguir una distribución exponencial de media 1 ($\beta=0$)
- Si $S^*(t)$ es el estimador K-M de estas observaciones entonces su $H^*(t)$ debe de seguir una línea de 45°.
- $H^*(t) = -\log(\exp(-t)) = t$

47

Residuos de Cox-Snell

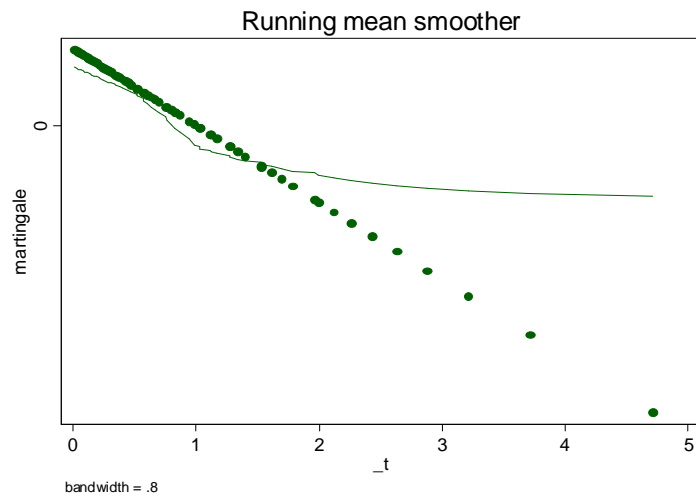
```
stcox grupo,mgale(mg)
predict cs,csnell
stset cs,f(mort)
sts generate km=s
gen H=-ln(km)
graph7 H cs cs ,c(11) s(...) xlab ylab
```



48

Residuos de Martingala

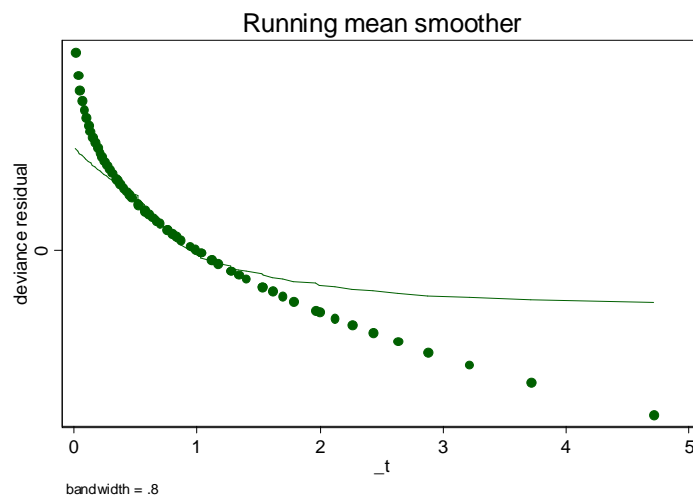
```
xi:stcox i.expcateg,mgale(mg)  
ksm mg _t, ylabel(0)
```



49

Residuos de Lejanía (Deviance)

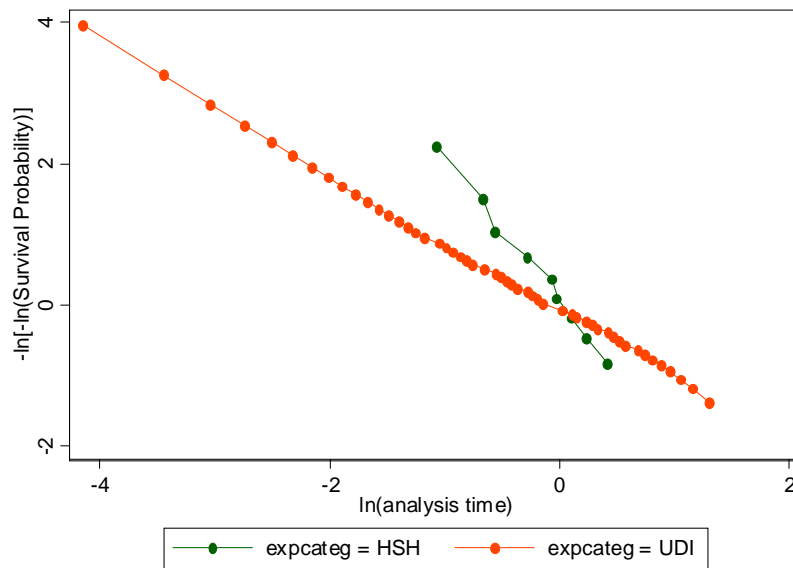
```
xi:stcox i.expcateg,mgale(mg)  
predict dev,deviance  
ksm dev _t, ylabel(0)
```



50

Gráfico log-log

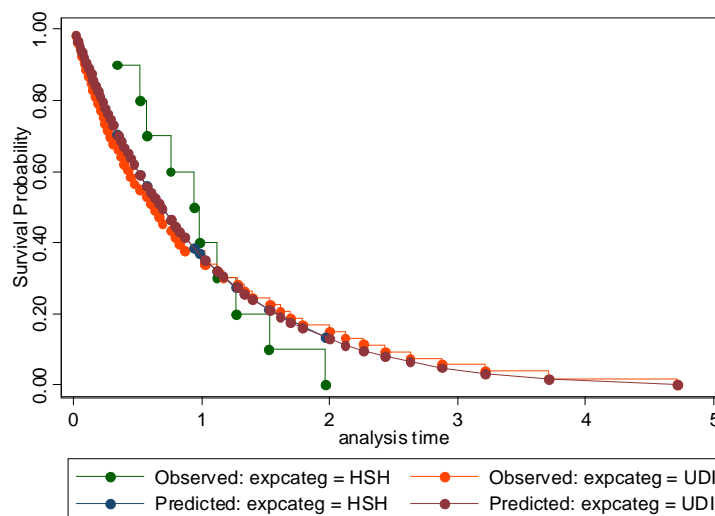
```
xi:stcox i.expcateg  
stphplot,by(expcateg)
```



51

Compara KM con predicción

```
xi:stcox i.expcateg  
stcoxkm,by(expcateg)
```



52

Test de contraste (Grambsch y Therneau)

```
xi:stcox i.expcateg,scale(sca*) schoenfeld(scho*)
stphtest , detail

stphtest , detail
```

Test of proportional-hazards assumption

Time: Time

	rho	chi2	df	Prob>chi2
_Iexpcateg_2	-0.18136	2.07	1	0.1504
global test		2.07	1	0.1504

Dibujar predicciones

```
stcurve, survival at1(_Iexpcateg_2=0) at2(_Iexpcateg_2=1)
```

