

1 Sex-Age-CalendarTime Patterns in population mortality rates in Denmark

Exercise 1: Use the same informal approach as earlier (OR – only if interested– a median polish), to fit a multiplicative model to the slightly larger dataset consisting of the 24 rates for all 3 periods i.e., to the data involving the 3 periods 1980-84, 2000-2004 and 2005-2007.

Yrs	Age	Female (F)	Male (M)
'80- '84	70-	R_F	$R_F \times M_M$
	75-	$R_F \times M_{75}$	$R_F \times M_{75} \times M_M$
	80-	$R_F \times M_{80}$	$R_F \times M_{80} \times M_M$
	85-	$R_F \times M_{85}$	$R_F \times M_{85} \times M_M$
'00- '04	70-	$R_F \times M_{20y}$	$R_F \times M_M \times M_{20y}$
	75-	$R_F \times M_{75} \times M_{20y}$	$R_F \times M_{75} \times M_M \times M_{20y}$
	80-	$R_F \times M_{80} \times M_{20y}$	$R_F \times M_{80} \times M_M \times M_{20y}$
	85-	$R_F \times M_{85} \times M_{20y}$	$R_F \times M_{85} \times M_M \times M_{20y}$
'05- '07	70-	$R_F \times M_{25y}$	$R_F \times M_M \times M_{25y}$
	75-	$R_F \times M_{75} \times M_{25y}$	$R_F \times M_{75} \times M_M \times M_{25y}$
	80-	$R_F \times M_{80} \times M_{25y}$	$R_F \times M_{80} \times M_M \times M_{25y}$
	85-	$R_F \times M_{85} \times M_{25y}$	$R_F \times M_{85} \times M_M \times M_{25y}$

R = rate. M = multiplier. The array called 'r' in the R code (which fits additive models to the rates and logs of the rates) can be used to calculate ratios.

...Year.....Age...Female...Male.....Total... Observed rates

1980-1984 70-74 0.02725 0.05213 0.03814
 1980-1984 75-79 0.04592 0.08235 0.06042
 1980-1984 80-84 0.08098 0.12163 0.09561
 1980-1984 85-89 0.13680 0.18202 0.15193

2000-2004 70-74 0.02666 0.03972 0.03261
 2000-2004 75-79 0.04179 0.06586 0.05189
 2000-2004 80-84 0.06923 0.10584 0.08279
 2000-2004 85-89 0.11970 0.16773 0.13480

2005-2007 70-74 0.02359 0.03468 0.02874
 2005-2007 75-79 0.03934 0.05815 0.04750
 2005-2007 80-84 0.06559 0.09622 0.07730
 2005-2007 85-89 0.11462 0.15808 0.12860

Age multipliers:

The rate in the (females 70-74, 1980-84) cell is 0.02725, while that in the cell one below it (75-79) is 0.04592, yielding an empirical rate ratio of 1.69 for the pure 75-79 vs 70-74 contrast. We can repeat the same 75-79 vs 70-74 contrast for each of the other 5 sex-calendar year combinations, to obtain in all six 75-79 vs 70-74 ratios:

Years	Age	Female (F)	Male (M)
1980-1984	70-74	1	1
	75-79	1.69	1.57
2000-2004	70-74	1	1
	75-79	1.58	1.66
2005-2007	70-74	1	1
	75-79	1.67	1.68

One way, without even using a calculator, to arrive at a best estimate of the M_{75} multiplier is to make the median, 1.66, of these 6 estimates.

Moving on to the the pure 80-84 versus 70-74 contrast, we obtain 6 rate ratio estimates: 2.97, 2.60, 2.33, 2.66, 2.78 and 2.77; their median is 2.72.

For the 85-89 versus 70-74 contrast, the median of the 6 estimates is 4.52.

These three multipliers can be used to derive multiplicative rate (i.e., insurance premium) increases for the higher age categories, using the rates in the 70-74 group as the reference or 'starter' or 'corner' category ('corner' is Clayton and Hills terminology in their chapter 22).

It seems that rates double about every 7 years or so. Note also that the estimated 10 year increase of 2.72 is virtually the same as 1.66^2 , so in fact we could use two 66% 5-year increases, 1 each per 5 years of age, and avoid having (to memorize/estimate) a separate multiplier for the 10 years of age increase. Note also that $1.66^3 = 4.57$ which is quite close to the fitted 4.52. So, in fact we could save having to memorize not just 1 but 2 multipliers, and simply say the rates in those ages 75-79, 80-84 and 85-89 are 1.66, 1.66^2 , and 1.66^3 times the rates in those aged 70-74.

Another way to say this is that the *logs* of the mortality rates are *linear* in *age*. This finding is not new: The actuary Benjamin Gompertz described this pattern as a Law of Mortality (that now bears his name) in a paper in 1825. And William Farr and Thomas R Edmonds, and Gompertz, used this smooth

functions relationship to save a lot of steps in the otherwise tedious life table calculations used in actuarial and population-life table analyses. When we come to formally fitting multiplicative rate (ie log linear) models for rates, the fact that the log rates seem to be close to linear over this age range means that we do not have to model age as a ‘categorical’ variable with 3 indicator variables (3 separate coefficients) but instead can be parsimonious (economical, even frugal) and use just 1 linear age term and its 1 associated regression coefficient.

Male multiplier:

The rate in the (females 70-74, 1980-84) cell is 0.02725, while that in the cell to the right of it (Males) is 0.05213, yielding an empirical rate ratio of 1.91 for the pure M vs F contrast. We can repeat the same M vs F contrast for each of the other 11 age-calendar year combinations, to obtain in all twelve M vs F ratios:

Yrs	Age	Female (F)	Male (M)
	70-74	1	1.91
'80-	75-79	1	1.79
'84	80-85	1	1.50
	85-90	1	1.33
	70-74	1	1.49
'00-	75-79	1	1.58
'04	80-84	1	1.53
	85-	1	1.40
	70-74	1	1.47
'05-	75-79	1	1.48
'07	80-84	1	1.47
	85-	1	1.38

The median of these 12 estimates is 1.48; one interpretation is that males should pay 48% higher life insurance premiums than females!

20-year multiplier: unchanged from in smaller dataset

The rate in the (females 70-74, 1980-84) cell is 0.02725, while that in the cell 4 cells below it (also females-70-74, but 20 years later) is 0.02666, yielding an empirical rate ratio of 0.98 for the pure ‘20 calendar years’ contrast. We can repeat the same contrast for each of the other 7 age-sex combinations, to obtain in all eight 2000-2004 vs 1980-1984 ratios:

Age	Female (F)	Male (M)
70-74	0.98	0.76
75-79	0.91	0.80
80-84	0.85	0.87
85-89	0.88	0.92

The median of these 8 estimates is 0.88 representing a reduction of 12% in mortality in the 20 years between 1980-1984 and 2000-2004.

25 (24?)-year multiplier:

The rate in the (females 70-74, 1980-84) cell is 0.02725, while that in the cell 8 cells below it (also females-70-74, but 24 years later) is 0.02359, yielding an empirical rate ratio of 0.87 for the pure ‘24 calendar years’ contrast. We can repeat the same contrast for each of the other 7 age-sex combinations, to obtain in all eight 2005-2007 vs 1980-1984 ratios:

Age	Female (F)	Male (M)
70-74	0.98	0.66
75-79	0.86	0.71
80-84	0.81	0.79
85-89	0.84	0.87

The median of these 8 estimates is 0.82 representing a reduction of 18% in mortality in the 24 years between 1980-1984 and 2005-2007.

corner term (a.k.a. the ‘intercept’:

Whereas all of the other estimates used a synthesis of several estimates, it is not immediately obvious whether we are forced to use the one observed value in the ‘corner’ cell as the best fitted value for that cell. But for now, let's use it as the corner estimate, so that we can write a master equation for all 24 rates

The equation is for the rate in any given age-group in a given gender in a given calendar period:

Rate =	0.02725	×1.66	×2.72	×4.52	×1.48	×0.88	×0.82	-3.505	-3.106
		if	if	if	if	if	if	-3.005	-2.606
		75-79	80-84	85-89	male	2000-04	2005-07	-2.510	-2.111
								-2.002	-1.603
log[Rate] =	-3.603	+0.509	+1.000	+1.509	+0.395	-0.136	-0.194		
		if	if	if	if	if	if		
		75-79	80-84	85-89	male	2000-04	2005-07	-3.633	-3.234
								-3.133	-2.734
log[Rate] =	β_0	$+\beta_{75'}$	$+\beta_{80'}$	$+\beta_{84'}$	$+\beta_M$	$+\beta_{20y'}$	$+\beta_{25y'}$	-2.637	-2.239
		×	×	×	×	×	×	-2.130	-1.731
		I_{75-79}	I_{80-84}	I_{85-89}	I_{male}	$I_{2000-04}$	$I_{2005-07}$		
								-3.739	-3.340
								-3.240	-2.841
								-2.744	-2.345
								-2.236	-1.838

where each ' I ' is a (0/1) indicator of the category in question.

By using both the 0 and 1 values of each I , this 7-parameter equation produces a fitted value for each of the $4 \times 2 \times 3 = 24$ cells.

You can also think of I_{75-79} , I_{80-84} , and I_{85-89} as 'radio buttons': at most 1 of them can be 'on' at the same time, since there are 4 age levels in all.

1.1 More formal fitting of 6 parameter values

It shouldn't have to be, in the model fitting above, that the intercept was forced to go through an observed value, when we know that that value (like each of the 15 others) is subject to sampling variation. A fitted regression line or curve that goes *between* the dots [as opposed to one that actually *joins* the (error-containing!) dots] recognizes the fact that none of the observed data-points is 'perfect.' Also the purpose of the line is as a 'line of means' or 'line of centres.'

One option to avoid the arbitrariness in fitting an intercept is to apply a **median polish** to the log-rates. You can look up this procedure on the web, and the c634 course website provides some code for carrying it out (It seems that the `medpolish` function in R just handles 2 dimensional arrays, whereas the homemade R function is designed for ≥ 2 dimensions).

The fitted values from the median polish of the $4 \times 2 \times 3$ array of log rates are given in the next column

Converting them back to rates, and scaling them all so that the corner is 1, we get the following fitted rate ratio model:-

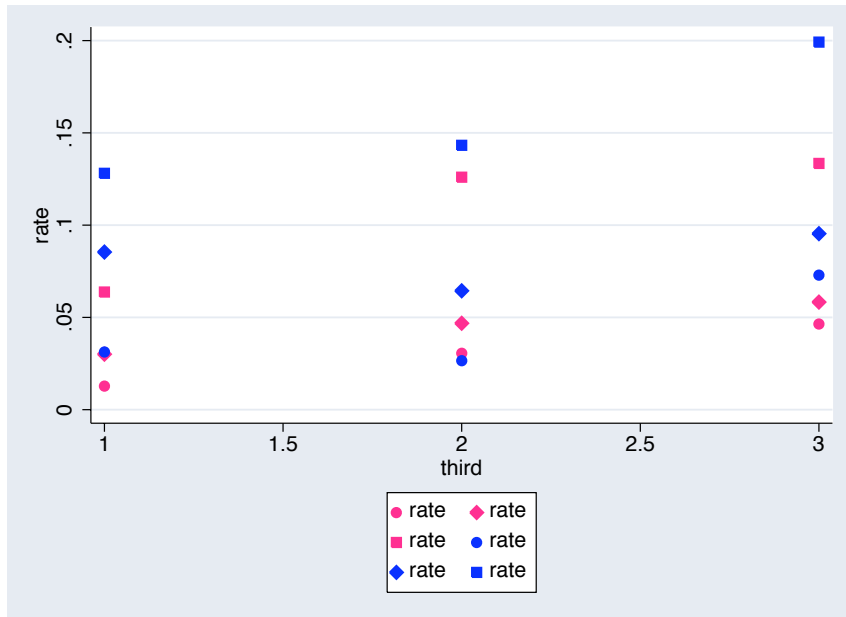
RateRatio =	1	×1.68	×2.71	×4.49	×1.49	×0.88	×0.79
		if	if	if	if	if	if
		75-79	80-84	85-89	male	2000-04	2005-07

2 Comparison of ≥ 2 Rates - via regression

Exercise 2 : data in Tables 1 and 2 in the Perceived-Age article.

- i. Within each of the 6 sex-age strata, there are has 3 rates – one for each ‘third’ of the perceived-age distribution. Plot these 18 rates on a single graph, with ‘third’ (1 2 3) on the horizontal axis, the rate on the vertical axis, and using different symbols for the 6 strata.¹

Stata graphics... JH isn't very good at annotating them.. better in R – see website for some R code.

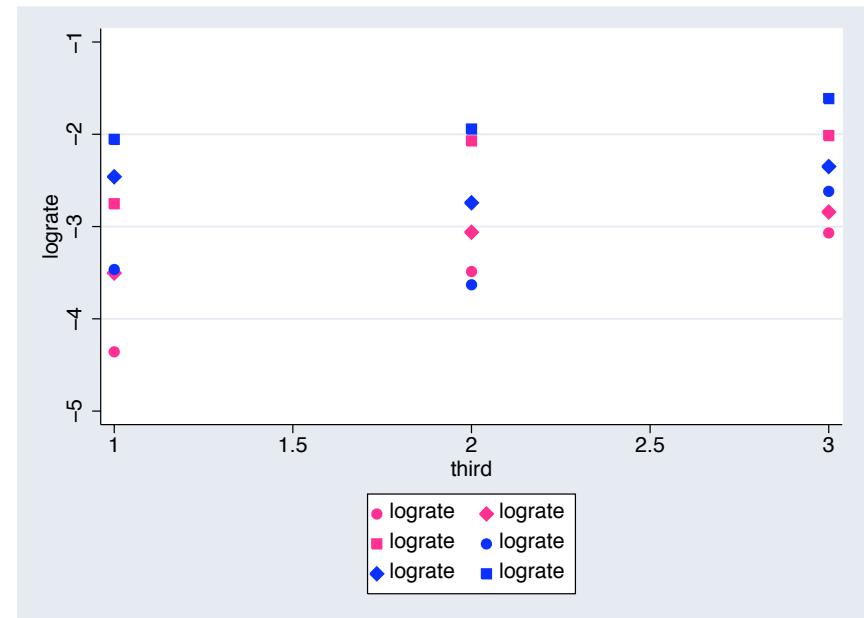


- ii. Re-plot these 18 rates on a new graph, but using a log scale for the rates.

```
scatter lograte third if ( male==0 & agecat==1), mcolor(pink) msymbol(circle) ///
|| scatter lograte third if ( male==0 & agecat==2), mcolor(pink) msymbol(diamond) ///
|| scatter lograte third if ( male==0 & agecat==3), mcolor(pink) msymbol(square) ///
|| scatter lograte third if ( male==1 & agecat==1), mcolor(blue) msymbol(circle) ///
|| scatter lograte third if ( male==1 & agecat==2), mcolor(blue) msymbol(diamond) ///
|| scatter lograte third if ( male==1 & agecat==3), mcolor(blue) msymbol(square)
```

There must be a simpler way, and one that allows better legends.

¹The rates resources on the c634 website has R code that can create the plots. Or you might wish to use Stata.



- iii. By eye, fit 6 parallel lines to the 18 (6 sets of) $\log(\text{rate})$'s. Log.rates are more parallel than rates themselves, which spread out with age.

Eyes will differ; JH saw a common slope of approx. 0.3 per third.

- iv. Using the multiple regression package of your choice, fit an additive model to the 18 $\log(\text{rate})$'s. Then convert it to a 'multiplicative rates' model. Ignore for the moment the fact that each log-rate is measured with a different precision.

Stata listing of data is on next page

Curious to know what if (first) just estimate the (crude) difference in log rates b/w men and women... Do not report as many decimals as Stata does.. JH didn't have time to figure out how to get Stata to report fewer!

```
. reg lograte male
```

Source	SS	df	MS	Number of obs =	18
Model	1.01893424	1	1.01893424	F(1, 16) =	2.08
Residual	7.84874154	16	.490546346	Prob > F =	0.1688
Total	8.86767578	17	.521627987	R-squared =	0.1149
				Adj R-squared =	0.0596
				Root MSE =	.70039

lograte	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
male	.4758464	.3301671	1.44	0.169	-.2240765 1.175769
_cons	-3.017957	.2334634	-12.93	0.000	-3.512878 -2.523037

Coefficient is 0.4758464, which translates to a crude rate ratio of $\exp(0.4758464) = 1.61$, not far from the age-adjusted one you calculated in an earlier exercise. This is not that surprising as the Coefficient is simply the difference b/w the mean of the 6 male log.rates and the mean of the 6 female log.rates. It is a form of standardization, with weights of 1/6 each.

Using 'thirds' as an interval variable...

```
. reg lograte agecat male third
```

Source	SS	df	MS	Number of obs =	18
Model	7.98897558	3	2.66299186	F(3, 14) =	42.43
Residual	.878700199	14	.0627643	Prob > F =	0.0000
Total	8.86767578	17	.521627987	R-squared =	0.9009
				Adj R-squared =	0.8797
				Root MSE =	.25053

lograte	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
agecat	.6817131	.0723212	9.43	0.000	.5265996 .8368267
male	.4758464	.1181001	4.03	0.001	.222547 .7291459
third	.3407403	.0723212	4.71	0.000	.1856268 .4958539
_cons	-5.062864	.220945	-22.91	0.000	-5.536744 -4.588984

log.rates are approx. .34 higher per third, so rates are approx. $\exp(.34) = 1.4$ higher per third.

Using 'thirds' as a categorical variable...

```
. reg lograte I_AgeC2 I_AgeC3 male I_Third2 I_Third3
```

Source	SS	df	MS	Number of obs =	18
Model	8.02502222	5	1.60500444	F(5, 12) =	22.86
Residual	.842653561	12	.07022113	Prob > F =	0.0000
Total	8.86767578	17	.521627987	R-squared =	0.9050
				Adj R-squared =	0.8654
				Root MSE =	.26499

lograte	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
I_AgeC2	.6117619	.1529936	4.00	0.002	.2784175 .9451063
I_AgeC3	1.363426	.1529936	8.91	0.000	1.030082 1.696771
male	.4758464	.1249188	3.81	0.002	.2036718 .748021
I_Third2	.2765649	.1529936	1.81	0.096	-.0567795 .6099094
I_Third3	.6814807	.1529936	4.45	0.001	.3481363 1.014825
_cons	-3.995702	.1529936	-26.12	0.000	-4.329046 -3.662357

log.rates are approx. $\exp(.27) = 1.3$ higher in 2nd third than first, and $\exp(.68) = 2$ times higher in 3rd third than first.

. list

	male	third	agecat	ndeaths	pyears	rate	lograte	I_AgeC1	I_AgeC2	I_AgeC3	I_Third1	I_Third2	I_Third3
1.	1	1	1	25	800.3	.0312383	-3.466111	1	0	0	1	0	0
2.	1	2	1	21	792.3	.0265051	-3.630418	1	0	0	0	1	0
3.	1	3	1	49	672	.0729167	-2.618438	1	0	0	0	0	1
4.	1	1	2	39	456.6	.0854139	-2.460246	0	1	0	1	0	0
5.	1	2	2	30	465.7	.0644192	-2.742344	0	1	0	0	1	0
6.	1	3	2	42	440.3	.0953895	-2.349787	0	1	0	0	0	1
7.	1	1	3	42	327.9	.1280878	-2.055039	0	0	1	1	0	0
8.	1	2	3	46	321.1	.1432575	-1.943111	0	0	1	0	1	0
9.	1	3	3	54	271.1	.1991885	-1.613504	0	0	1	0	0	1
10.	0	1	1	10	781.2	.0128008	-4.358246	1	0	0	1	0	0
11.	0	2	1	23	752.5	.0305648	-3.487907	1	0	0	0	1	0
12.	0	3	1	33	710.5	.0464462	-3.069461	1	0	0	0	0	1
13.	0	1	2	18	598.5	.0300752	-3.504055	0	1	0	1	0	0
14.	0	2	2	27	576.7	.0468181	-3.061485	0	1	0	0	1	0
15.	0	3	2	31	531.7	.0583036	-2.842092	0	1	0	0	0	1
16.	0	1	3	42	658.7	.063762	-2.752599	0	0	1	1	0	0
17.	0	2	3	74	587.4	.1259789	-2.071641	0	0	1	0	1	0
18.	0	3	3	69	517.1	.1334365	-2.01413	0	0	1	0	0	1

Structured 2-D or 3-D datasets where the (even approx.) additivity of log rates (multiplicative pattern of rates) does not hold.

The handout “Survival (or Cumulative Incidence*) functions, v. 2010.01.21” has several examples.

In the Uganda trial of male circumcision, the reduction in the rate of HIV transmission from mo. 12 onwards is greater than that in mo. 1-6, & 7-12.

In the comparisons of mortality rates in MI patients, the increased rate in those admitted on weekends versus weekdays is limited to the first 3-4 days. After that, the mortality rates are virtually identical.

In the European RCT of prostate cancer screening, prostate cancer mortality in the screened arm was virtually identical to that in the control arm for the first 7 years, but substantially less after that. The 20% reduction reported in the abstract is an inaccurate figure. The analysis that was used in this screening study, and in many other cancer screening studies to date, is the statistical equivalent of measuring the long-term steady state reduction in a patient’s LDL cholesterol as a result of statin treatment by comparing the average pre-treatment LDL level with a new average based mainly on the levels in the early (post-initiation-of-statins) period before the LDL reaches a new steady state.