



Referanser (forts.)

Rässler, S., Rubin, D. B., & Zell, E. R. 2008, "Incomplete Data in Epidemiology and Medical Statistics," in Epidemiology and Medical Statistics (Handbook of Statistics vol 27)". C. R. Rao, J. P. Miller, & D. C. Rao, eds., Elsevier, pp. 569-601.

Carpenter, J. R. & Kenward, M. G. (2007) Missing data in randomised controlled trials - a practical guide. 186 pages. http://www.pcpoh.bham.ac.uk/publichealth/methodology/proj ects/RM03\_JH17\_MK.shtml

Little, R J A, Rubin, D B: (2002) *Statistical Analysis with Missing data*. 2nd ed. Wiley.

Schafer, J. L. (1997) Analysis of incomplete mul Carle dae U Chapman & Hall, London.

# Manglende data:

- "Hull" i datamatrisen, som ideelt sett burde være komplett
- Vanligvis data som man hadde til hensikt å registrere, men som av en eller annen grunn ikke ble registrert.
- Det finnes en meningsfylt dataverdi som ikke er registrert.

NTNU



"Less than optimum strategies for missing values can produce biased estimated, distorted statistical power, and invalid conclusions. After reviewing traditional approaches (listwise, pairwise, and mean substitution), selected alternatives are covered including single imputation, multiple imputation, and full information maximum likelihood estimation. ... When missing values cannot be avoided, multiple imputation and full information methods offer substantial improvements over traditional approaches."

NTNU









# Missing data problem: 987 measurements (7 time points x47 patients x 3 substances) Missing data for 7 of 987 measurements This was 4 of the 47 patients! Repeated measurements ANOVA (as used in this study) requires complete data







45			
10			
Variable	n	% missing	
Systolic BP	64708	1,3	
Diastolic BP	64708	1,3	
Cholesterol	65158	0,7	
HDL-Cholesterol	65155	0,7	
GLUCOSE	65158	0,7	
Triglycerides	65158	0,7	
Creatinine	65158	0,7	
eGFR <sup>1)</sup>	65158	0,7	
ACR <sup>2)</sup>	9703	85,2	
<sup>1)</sup> estimated glomeru <sup>2)</sup> Albumin creatinin Not requested (Mi Requested, but no	ular filtratior ratio (from ssing by de ot deliverd: :	n rate urine sample) ssign): 82,8 % 2,5%	D NTNU Det skagende universitet
www.ntnu.no			























European or cancer QLQ	ganization for re -C30 (EORTC (	esearch a QLQ-C30)	ind trea	atment o	of	
		Not at all	A little	Quite a bit	Very much	]::
21. Did you f	eel tense?	1.	2	з	4	
22. Did you v	vorry?	1	2	з	4	1.0
23. Did you f	eel irritable?	1	21	3	4	
24. Did you fe	eel depressed?	1	2√	3	4	
Figure 15.3 The em	otional functioning s	cale of the E	ORTC Q	LQ-C30		
and the second second				Det s	TNU	sitet
www.ntnu.no						























45	46
ML estimation	ML Estimation
Some theory	When data are MAR, the observed data-likelihood is
$P(Y_{obs}; \theta) = \int P(Y_{com}; \theta) dY_{mis}$	$L(\theta; Y_{obs}) = P(Y_{obs}; \theta)$
is a correct - probability distribution if MCAR - likelihood if MAR or missing values are out of scope	The ML estimate $\hat{\theta}$ that maximises the likelihood.
$P(Y_{obs}, R; \theta, \xi) = \int P(Y_{com}; \theta) P(R \mid Y_{com}; \xi) dY_{mis}$	The log likelihood may be easier to calculate:
is a correct likelihood if MNAR (difficult)	$l(\theta; Y_{obs}) = \log L(\theta; Y_{obs})$
D NTNU Det skapende universitet	NTNU     Det skapende universitet
www.ntnu.ne	www.ntnu.ne



















Hallan & al (2009) Statistical analyses were performed using Stata 10.0 (Stata Corp., TX, U.S.A.). In general, there were few missing data (<2% for most variables, see Table 1), but data on ACR were, by study design, available only in a subgroup. Multiple imputation is now considered the standard method for handling this type of data,(Clark & Altman 2003;Donders et al. 2006;Rassler et al. 2008;van Buuren et al 1999) whereas complete case analysis would yield too imprecise as well as biased results. The multiple imputation technique estimates the mean and uncertainty of the missing data using all information from the actually observed data in a proper way. In this way, unbiased estimates with the correct standard deviation and p-values are calculated.(Rassler et al, 2008).	Continued: For most nor missing com urine sample at random, ti analyses we procedures f and not used variables,(va were include the time vari 1999a) Regr sex and both these two inf We used m= accuracy.(No
A A A A A A A A A A A A A A A A A A A	and the second sec

n-diabetic non-hypertensive subjects data were pletely at random, and for those not returning es as requested data were assumed to be missing hus meeting the assumptions for the method. The re carried out in the "ice" and "micombine" for Stata, (Royston 2005) ACR was log-transformed d as predictor in the imputation of other missing an Buuren et al 1999b) study outcome variables ed in the imputation model, (Moons et al. 2006) and able was log-transformed (van Buuren, et al ression modelling revealed interactions between blood pressure and diabetes mellitus. Hence, teractions were included in the imputation model. =20 imputations to achieve maximum ewgard & Haukoos 2007) NTNU

## MI implementation in Stata

- · ice each variable with missing data is predicted from the other variables using the appropriate imputation regresson models (Linear, Binary logistic, or Ordinal logistic) and creates m imputed (complete) data sets
- micombine analyses each imputed data set using the relevant revant analysis model (for example Cox proportional hazards regression), and combines the results using Rubin's rules

NTNU

mim (2008) replaces micombine. Several new features.

# Example Hallan et al 2009. Implementation in Stata (ice) Categorical variables must be coded 0,...k-1. For example female is coded 0 and 1 $\,$ Continuous variables are assumed normally distributed. Used In(ACR) instead of ACR. Do not use a predictor with more than 50% missing. (Hence In(ACR) used only as dependent variable) Include outcome variable as predictor. Here: follow-up time and event CKD. Use log transformed time variable as predictor (outcome variable in the Cox analysis model) Do not impute outcome if missing! Use an imputation model at least as rich as the analysis model. We included the interactions sex\*bp and sex\*diabetes. Used a high number of imputations (m=20) due to high proportion missing.

NTNU



		Completed the physical-pe	erformance evaluation	
	Yes		No	
Dead by 12/31/1991	No	Yes	No	Yes
n	1527	264	416	134
Age, median (IQR)*	77 (74-81)	78 (74-84)	78 (74-83)	85 (78-90
Male, %	32.0	45.1	30.3	47.0
Physical performance median (IQR)	8 (5-10)	6 (2-8)	-	-
Self-assessed health median (IQR)	2 (2-3)	3 (2-3)	2 (2-3)	3 (2-3)
notes: These are the 2341 residents of East Bost	on who participated in the 6 noie, Gluno et al. [35]). They	I-year follow-up evaluation of ranged in age from 71 to 10	the EPESE, or Established I 73 years at that time. Particip	Populations for ants were asked to ation of physical
Epolemologic Studies of the Eddery (see, for easy rank the'n heah relative to others of the's age as 1 performance. This was based on brief stats of bala corea indicing better function (Guralnik et al. [36 "Interquantile range.	excellent; 2, good; 3, fair; o ance, gait, strength, and end [).	r 4, poer. Those who were a lurance, and summarized in a	ore, also had objective evalu in overall score ranging from	0 to 12, with high
Epoemoogic Studies of the Edely (see, for easy and their health relative to other of their age as 1 performance. This was based on brief tests of bala cores indicing better function (Guusinia et al. [36 friterquartile ranges.	excellent; 2, good 3, fair, anon, gait, strength, and end D.	r 4, poor. Those who were a lurance, and summarized in a	ow, also had cojective evalu n overall score ranging from	0 to 12, with higher $\Gamma NU$



A complete case analysis (n=7,832 and 43 kidney failure

Continued:

	Imputatio	n number				Total, by Rubin's rules	p- value	FMI
	1	2	3	4	5			1
Age,	0.0707	0.0705	0.0701	0.0701	0.0707	0.0704	< 0.001	0.002
years	(0.0067)	(0.0067)	(0.0067)	(0.0067)	(0.0067)	(0.0067)		
Female	-0.612	-0.580	-0.570	-0.582	-0.589	-0.587	0.002	0.008
sex	(0.189)	(0.190)	(0.190)	(0.190)	(0.190)	(0.191)		
ACR	0.0276	0.0282	0.0285	0.0283	0.0280	0.0281	< 0.001	0.082
	(0.0013)	(0.0013)	(0.0013)	(0.0013)	(0.0013)	(0.0014)		

5 of the imputations in Hellen et al 2000)





		Analytic n	nethod	
	Complete case	No physical performance	Multiple imputation	Indicator method
n (deaths)	1791 (264)	2341 (398)	2341 (398)	2341 (398)
Age	0.033 (0.013)	0.088 (0.009)	0.057 (0.011)	0.068 (0.01)
Male	0.92 (0.15)	0.82 (0.12)	1.00 (0.14)	0.92 (0.12)
Self-assessed health	0.38 (0.095)	0.60 (0.073)	0.39 (0.087)	0.46 (0.076)
Physical performance	-0.14 (0.023)	-	-0.15 (0.024)	-0.12(0.022)
Intercept	-4.76 (1.17)	-10.41 (0.79)	-6.60 (1.04)	-7.42 (0.92)
Indicator of missing performance				-0.47 (0.13)

		5 drawn values of parameters				
Variable	Estimate (SE)	1	2	3	4	5
Age (per year)	-0.22 (.013)	-0.23	-0.23	-0.23	-0.21	-0.22
Male gender	1.18 (0.15)	1.21	1.08	1.29	1.25	1.22
Self-assessed health	-1.38 (0.089)	-1.33	-1.43	-1.24	-1.43	-1.35
Dead by 1991	-1.30 (0.20)	-1.12	-1.47	-1.37	-1.60	-1.05
Intercept	27.2 (1.1)	27.6	27.7	27.5	26.3	27.2
Overall R <sup>2</sup> 0.30; Root N	ASE 2.88					











"The days that journals tolerate the absence of analysis of missing values and the use of traditional approaches to missing values should be numbered. ... In general, multiple imputation and the approaches available in structural equation modelling software are the best that are currently available."

NTNU

