GRELL & ENCOR 2025 Scientific Meeting Porto, Portugal



Privacy-preserving self-service linkage of cancer registry and socioeconomic data in a safe processing environment

Siri Larønningen

Cancer Registry of Norway, Norwegian Institute of Public Health



Background

Problems:

- Time consuming applications
- Lacking metadata
- Lacking possibility to explore data before the application ("cohort-explorer")
- Physical sharing of pseudonymized data
- Limited control over the data after it has been shared
- → It isn't always necessary to see the actual individual data

One possible solution: self-service in a secure environment

microdata.no – data access without application

microdata.no provides instant, online access to <u>large amounts of detailed and</u> <u>mergeable microdata</u> without any form of application.

The service is open to employees and students at universities and colleges, approved research institutions, ministries and directorates.

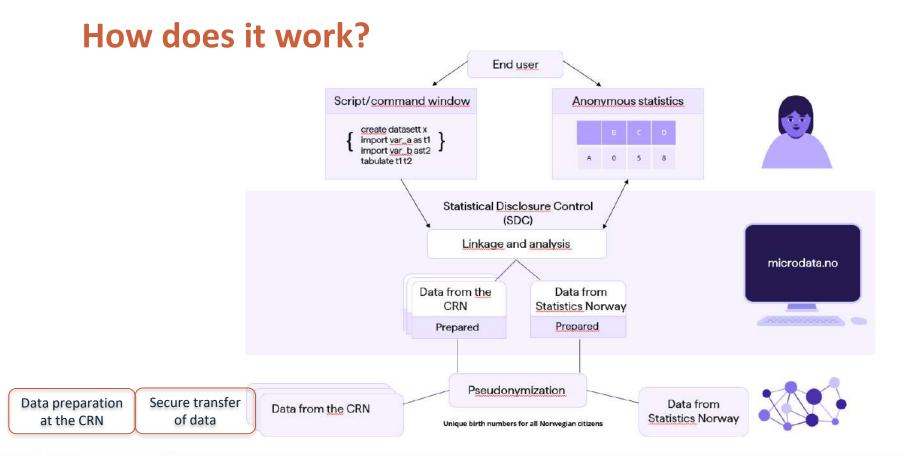
- No applications
- Instant access
- Time series from 1964
- Export function. Create a dataset and apply for it.
- · Self-service. The institutions register their users themselves.

A solution supporting the FAIR-principles!

Microdata.no

- Has been in use with data from Statistics Norway since spring 2018
- Currently more than 700 users from 75 institutions
- 350 publications have used statistics made in microdata.no (level 2 publications, PhDs, reports, reviews, master thesis etc.)
- 665 variables from Statistics Norway is currently available

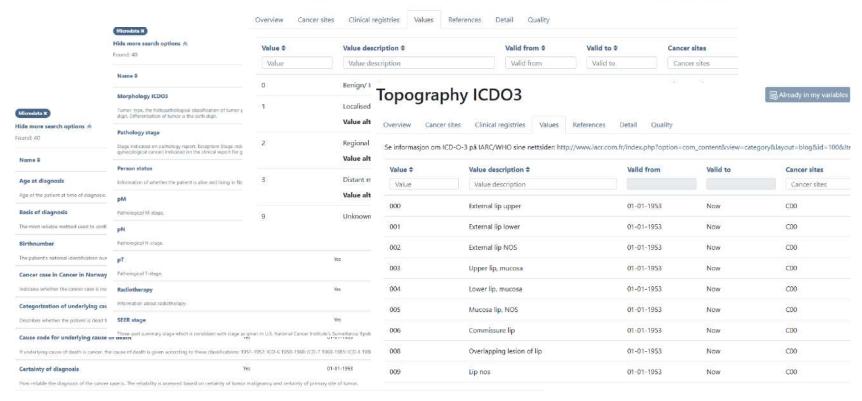
 AND: From Autumn 2025, variables from the Cancer Registry of Norway and from the Norwegian Patient Registry will also be available



Metadata

SEER stage





1	//Connect to KRG-databank (DRAFT = available for datatesters and dataadmins	-2//2						
	only)	Ì		see		erstadiur	n	
2	require no.fhi.kreftregisteret:DRAFT as kreg		0,,	1	2	3	9	Total
	//Connect to Statistics Norway-databank	~ 11		20121	10.0	1912 E-101	72272	
	require no.ssb.fdb:34 as fdb	0.5	1 - Lang høyere. Mor eller far eller begge har utdanning på nivå 7 el. 8	53.74	21.96	12.31	10.9	100
5	//Cancer incident dataset	uno	2 - Kort høyere. Mor eller far eller begge har utdanning på nivå 6	54.85	22.75	11.2	11.32	100
	create-dataset cancer_incident	kgu						
8	Greate datases cancer_and acre	bac	3 - Videregående. Mor eller far eller begge har utdanning på nivå 3, 4 eller 5	51.39	23.53	13.13	11.94	100
	//Importing date of diagnosis	cial	4 - Grunnskole. Mor eller far eller begge har utdanning på nivå 0, 1 eller 2	46.31	24.35	16.18	13.19	100
10.50	// import kreg/S_DIAGNOSEDATO as diagdato	200	and the state of t					
11	//Connect to KRG-databank (DRAFT = available for datatesters and dataadmins only) require no.fhi.kreftregisteret:DRAFT as kreg //Connect to Statistics Norway-databank require no.ssb.fdb:34 as fdb //Cancer incident dataset create-dataset cancer_incident //Importing date of diagnosis // import kreg/S_DIAGNOSEDATO as diagdato import kreg/S_DIAGNOSEDATO as diagdato		9 - Uoppgitt. Begge foreldrene har uoppgitt utdanning.	42.22	20	15.56	17.78	100
	//Using the year()-function to extract year from epoch days in variable		Total	10 83	23.68	14 13	12 27	100
	// generate diagnoseyear = year(diag ata)		CV-SMI	73.03	25.00	13110	72.07	100
15	generate diagnoseyear = year(diagnato)	Towns and the						
16	A	cancer_:	incident*barchart (percent) seerstadium, over(social_background) stack horizontal	L				
5,230	//Frequency table of ancouncidents pr year	v 1-	Lang høyere. Mor eller far ell		1 2			
18	tabulate diagnoleyea	uno 2	- Kort høyere. Mor eller far ell		3			
0.000	//Keeping only incidents from 2022 for now	200	- Videregående. Mor eller far	1	9			
21		-0						
22	(6)	OCI	- Grunnskole, Mor eller far ell					
23	Amport S_SEERSTADIUM (Summarisk tredeling av stadium i samsvar med U.S.	vi 9.	- Uoppgitt, Begge foreldrene					
	National Cancer Institute's SEER Program)		0 30 30 80 50 50 30 80 90 30	20 320				
	import kreg/S_SEERSTADIUM as seerstadium		percent	0				
25	//Import av "krefttilfellets fødselsnummer (pseudonymt)"							
	import kreg/P_FODSELSNR as kreg_fnr		<pre>incident* drop if inlist(highest_education_1_digit22, "0", "1")</pre>					
28	Import in Egy _ observant as in eg_ m	176 er	nheter ble fjernet fra datasettet.					
29	//Person dataset with data from SSB	cancer_	<pre>incident> barchart (percent) seerstadium, over(highest_education_1_digit22) stack</pre>	horizon	tal			
165	create-dataset person	2	2 - Ungdomsskoleutdanning		1			
1000	// *** RESIDENTS ONLY ***	git 2	Videregående, grunnutdanni		2			
	//Importing the STATUSKODE variable that lets us	9	Security of Page 200 to Security Securi		9			
	//keep only the population that lives in the country on	L C	4 - Videregående, avsluttende					
100	//a given date import fdb/BEFOLKNING_STATUSKODE 2022-01-01 as regstat22	catio	5 - Påbygg videregående					
	keep if regstat22 == '1'	npa	6 - Bachelor					
37	THE DEVICE WEST CONTROL OF THE STATE OF THE	est	7 - Master-					
38	// *** EDUCATION ***	high						
39	//Import highest education by 2022	E	8 - Forskerutdanning					
40	<pre>import fdb/NUDB_BU 2022-01-01 as highest_education_22</pre>		0 30 30 30 50 50 30 30 30 30	30 PS				
41			II A THE STATE OF	ಪ ಚೆ				



Security in microdata.no

The Five Safes framework:

Safe data: data is treated to protect any confidentiality concerns

Safe projects: research projects are approved by data owners for the public good.

Safe people: researchers/users are trained and authorised to use data safely.

Safe settings: the environment prevents unauthorised use.

Safe outputs: screened and approved outputs that are non-disclosive.

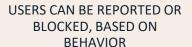
Statistical Disclosure Control (SDC)





Misuse detection







SYSTEM ADMINISTRATORS
CAN CHECK THE
REPORTED/BLOCKED SCRIPT
AND DECIDE ON ACTIONS



ACTIONS CAN BE WHITE-LISTED OR BLACK-LISTED



UTILITY VS. PRIVACY DISCUSSIONS



Pros and cons of microdata.no

PROs

- Instant accessibility if authorised user
- Linkage of health data and socioeconomic data
- Cohort-explorer
- Scripts can be shared and re-used (for instance with an application or individual data)
- Privacy-preserving

CONs

- Need to know how to write scripts and analyze data
- Not accessible for everyone (for instance drug companies, private companies etc)
- Limitations due to the SDC-layer
- Not suitable for complex analysis where direct access to data is crucial

Big thanks to the super-cool purple team from Statistics Norway, the Norwegian Agency for Shared Services in Education and Research, and the Cancer Registry of Norway

Statistics Norway:

Rune Gløersen, Vidar N. Klungre

Norwegian Agency for Shared Services in Education and Research: Ørnulf Risnes, Sigbjørn Revheim, Eirik Alvær, Eirik B. Stavestrand

Cancer Registry of Norway:

Jan F. Nygård, Sigrid Leithe, Narasimha Raghavan, Gintaras Pikelis, Siri Larønningen

