

Oppgavesett for innførende kurs i microdata.no - grunnleggende

Mål:

Bli kjent med de grunnleggende deler av microdata.no, med fokus på datasettoppbygging og deskriptiv statistikk. Til slutt prøver vi oss på en enkel regresjonsanalyse. Dette oppgavesettet passer for deg som ikke har programmeringserfaring.

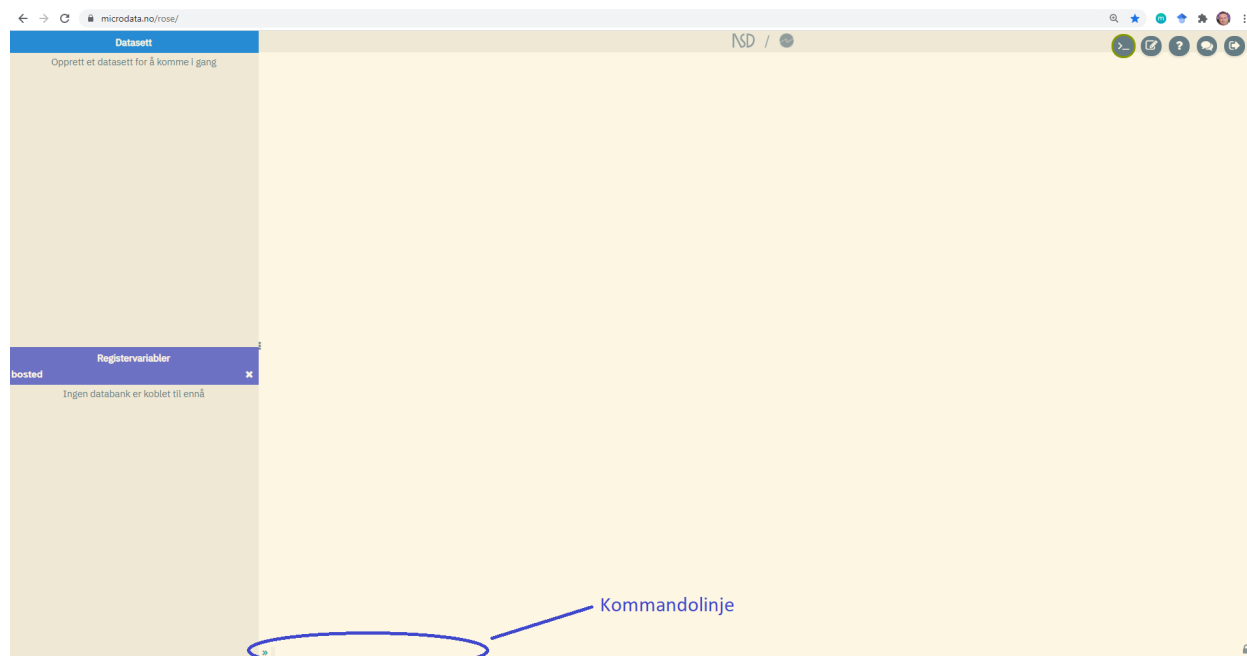
Innlogging/oppstart av microdata.no:

1. Skriv inn webadressen “microdata.no” i url-feltet i din nettleser. Bruk Chrome eller Firefox, evt. andre nettlesere bortsett fra Internet Explorer.
2. Hold musepeker over “Logg inn”-knappen og velg “Analyseverktøyet”
3. Velg innloggingsalternativ og logg inn på samme måte som ved en nettbank - vanligste er BankID på mobil (ha mobil tilgjengelig)

The screenshot shows the homepage of microdata.no. A blue circle highlights the address bar containing 'microdata.no', with a blue arrow pointing to it from below. Another blue circle highlights the 'LOGG INN' button in the top right navigation bar, with a blue arrow pointing to it from the right. The main content area features the heading 'microdata.no – registerdata uten å søke' and a list of features: 'Ingen søknader', 'Umiddelbar tilgang', 'Tidsserier fra 1964', 'Eksportfunksjon. Lag et datasett og søk om å få det utlevert.', and 'Selvbetjent. Institusjonene melder selv inn sine brukere.' Below this are buttons for 'BIBLOGRAFI' and 'VARIABLEOVERSIKT'. At the bottom, there are logos for Sikt, Kunnskapssektorens tjenesteleverandører, and Statistisk sentralbyrå (Statistics Norway).

Vi starter med å arbeide i analyseområdet i microdata.no. Dette er det området som møter deg når du har logget inn. Merk at før du har begynt å jobbe, er innholdet i vinduene tomme.

Helt nederst finner du kommandolinjen. Det er her en legger inn kommandoene en ønsker å eksekvere. Oppgave 7 i oppgavesettet går ut på å lagre alt arbeidet ditt som et skript. Dette gjør det mulig å lagre, systematisere, redigere og kjøre sett med kommandoer automatisk.



Nyttige tips:

Kommandoen `help` viser en oversikt over tilgjengelige kommandoer. `help` i kombinasjon med kommandonavn viser dessuten nyttig informasjon om den aktuelle kommando.

Tastekombinasjonen `<ctrl> + <z>` kan brukes til å angre og gå ett steg tilbake. Så om du gjør en feil, f.eks. koder en variabel på feil måte, kan du gjøre om på dette ved å gå ett trinn tilbake og kode på nytt. En kan i prinsippet klikke `<ctrl> + <z>` mange ganger og angre alt en har gjort.

Tasten "pil opp" på tastaturet gjør det mulig å hente frem tidligere brukte kommandoer. Disse kan så justeres på og kjøres som nye kommandoer.

For å se hvilket format verdiene til en variabel har, klikk på variabelnavnet i variabellisten for datasettet ditt i øverst til venstre i arbeidsområdet. Der finner du også kode-oversikt og enkel statistikk for den aktuelle variabelen. Det samme kan du gjøre i variabellisten nede til venstre. Denne inneholder alle tilgjengelige variabler i databasen (krever at du først bruker kommandoen `require` for å koble deg til databasen).

```
Variabler kan slettes eller døpes om etter behov. En bruker da hhv. kommandoene drop <variabel>  
og rename <navn_gml> <navn_ny>
```

I dette oppgavesettet skal vi se litt på hvordan lønnsinntekt har utviklet seg de siste årene, fordelt på ulike demografiske grupper av befolkningen.

Det fokuseres på å bygge opp datasett for analyser og å lage informativ deskriptiv statistikk. Til slutt demonstrerer vi hvordan vi avdekker statistiske sammenhenger mellom variabler, og kjører en enkel regresjonsanalyse for å finne årsakssammenhenger knyttet til lønnsinntekt.

Fasit finner du på de siste sidene av oppgavesettet.

1) Koble deg til databasen med variabler

Start alltid med å koble deg til databasen med variabler som du i de neste trinnene vil importere til ditt lokale datasett.

Per i dag er det kun én tilgjengelig database: SSBs database. Det vil etterhvert bli mulig å koble seg opp mot databaser som inneholder data også fra andre dataeiere.

Bruk følgende kommando i kommandolinjen for å koble deg til SSBs database:

```
require no.ssb.fdb:14 as ds
```

(Siste del av syntaxen, "as ds" brukes til å lage et kortnavn/alias på databasen. Her står du fritt til å velge andre navn om du ønsker.)

2) Opprett et tomt datasett

Neste trinn er å opprette et lokalt datasett som du skal fylle med variabler fra SSBs database. Bruk følgende kommando, og velg et passende navn på datasettet (erstatt <datasettnavn> med et eget navn):

```
create-dataset <datasettnavn>
```

Datasettet ditt vil først være tomt. Når du legger til variabler med kommandoen `import`, vil du se oppe i venstre vindu at datasettet fylles opp med stadig flere variabler, for hver import-kommando du eksekverer.

3) Definer populasjon

Etter å ha opprettet et datasett, er det på tide å definere populasjonen vi skal studere. Dette kan også gjøres senere, men det anbefales å gjøre dette så tidlig som mulig siden arbeidet med datatilrettelegging og analyser da blir mindre ressurs-/tidkrevende (jo mindre datasett/populasjon, jo raskere blir responsen).

Populasjoner defineres gjennom å:

- Hente variabler som brukes som populasjonskriterier
- Droppe enheter (individer) som faller utenfor kriteriene

Kommandoen `import` brukes til å hente variabler inn i ditt datasett fra databasen. Variabelopplysningene kobles automatisk opp mot respektive enheter (individer) i din populasjon vha. et innebygget identifikasjonsnummer, og dataene organiseres på tverrsnittsformat (én record per enhet):

```
import <databasealias>/<variabelnavn> <YYYY-MM-DD> [as <alias>]
```

<databasealias> erstattes med kortnavnet du har laget på SSBs database ("ds").
<variabelnavn> erstattes med navnet på variabelen du skal importere (husk store bokstaver), og *<YYYY-MM-DD>* med måletidspunktet, f.eks. 2019-12-31. Det er hensiktsmessig å bruke egne navn på variabler som importeres. Dette gjøres ved å legge inn "as" pluss et valgfritt navn til slutt i uttrykket.

Ved import av variabler med faste opplysninger (kjønn, fødeland, fødselsdato etc), skal ikke måletidspunkt angis.

Eksempel 1: `import ds/NUDB_BU 2019-07-31 as utd`

Eksempel 2: `import ds/BEFOLKNING_FOEDSELS_AAR_MND as faarmnd`

- a) Det er lurt å starte med en universell variabel, siden systemet benytter en left-join-logikk: Første variabel definerer din populasjon (som du senere kan trimme ned til ønsket størrelse).

Start med å importere variabelen `BEFOLKNING_FOEDSELS_AAR_MND` (fødselsdato på det numeriske formatet `YYYYMM`) til datasettet ditt. Bruk kommandoen `import`. Dette er en fast opplysning, så du skal ikke angi måletidspunkt.

Etter at import-kommandoen er kjørt, kan du se i vinduet øverst til venstre at datasettet består av fødselsdato-variabelen pluss en individ-nøkkelvariabel. Ved

å klikke på variabelnavnet dukker det opp en boks med metadata for denne variabelen.

Variabler med symbolet “123” foran har numerisk format, mens “abc” betyr at variabelen har tekstformat (alfanumerisk).

I samme boks vil du også se antallet records i datasettet ditt (bestemmes av populasjonen til den første importerte variabelen).

Du har nå laget et datasett som består av hele den norske populasjonen. Gjør deg kjent med populasjonen din ved å studere alderssammensetningen. Til dette trenger du først å lage en ny aldersvariabel basert på variabelen du har importert. Bruk kommandoen `generate` og mål alder per 2020. Variabel-aliaset til fødselsdato-variabelen som du importerte, brukes som input.

Forslag:

```
generate alder = 2020 - int(faarmnd/100)
```

(Det siste leddet i formelen trekker ut de fire første sifrene fra en 6-sifret fødselsdato, altså årstallet. Vi ønsker kun heltallet fra divisjonen. Derfor bruker vi funksjonen `int()`. Differansen mellom det oppgitte årstallet 2020 og fødselsåret blir da alder målt per 2020.)

b) Studer aldersfordelingen. Bruk kommandoen `histogram`:

```
histogram alder, discrete
```

Opsjonen `discrete` sørger for at alle alderskohorter får hver sin søyle i søylediagrammet. Dette anbefales når en skal bruke histogram til å se på diskrete verdier (og ikke intervaller av kontinuerlige verdier).

Legg merke til at aldersspennet strekker seg langt over 100. Grunnen til dette er at populasjonen din ikke er skikkelig filtrert. Datasettet ditt inneholder nemlig alle personer i databasen som noensinne har levd, inkludert døde personer.

c) Du trenger en variabel som angir om personer er bosatt, død eller emigrert for å fjerne personer du ikke ønsker å inkludere i din analyse. Til dette bruker du variabelen `ds/BEFOLKNING_STATUSKODE`. Denne variabelen er alfanumerisk og har verdiene ‘1’, ‘3’ og ‘5’ som betyr hhv. bosatt, utvandret og død. Merk også at variabelen er en såkalt tverrsnittsv variabel som bare har verdier for en fast dato

per år. Det varierer når de ulike tverrsnittsvariablene måles, men i dette tilfellet er det 1/1 som er det faste måletidspunktet (i kommandolinjen foreslås det gyldige datoer i autocomplete-menyen som dukker opp ved inntasting av kommandoen, og du kan dessuten sjekke gyldige måletidspunkter ved å klikke på variabelen i database-listen nederst til venstre, evt. i variabeloversikten du finner via innloggingssiden).

Importer først variabelen og bruk måletidspunktet 2020-01-01. Bruk deretter kommandoen `keep if` for å ta vare på kun dem som var bosatt i Norge på denne datoen (du fjerner altså personer som døde før 2020-01-01 eller som var utvandret på denne datoen).

Forslag:

```
import ds/BEFOLKNING_STATUSKODE 2020-01-01 as regstat  
keep if regstat == '1'
```

Bruk histogram-kommandoen på nytt for å se hvordan aldersfordelingen nå ser ut (husk opsjonen `discrete`).

Du kan også bruke kommandoen `summarize` for å hente ut oppsummerende statistikk for alder (snitt, standardavvik, median etc): `summarize alder`

- d) I dette oppgavesettet vil vi studere lønnsinntekt, og da kan det være greit å avgrense populasjonen ytterligere til personer i alderen 31-49 (da er de fleste ferdige med studier og er kommet i jobb).

Bruk kommandoen `keep if` eller `drop if` i kombinasjon med standard IF-betingelser til å trimme ned datasettet til å inneholde aldersgruppen 31-49.

Forslag:

```
keep if alder > 30 & alder < 50
```

Se på antallet enheter i datasettet ditt. Det har nå blitt redusert fra 10.626.204 til 1.378.664 individer.

Ved logiske if-uttrykk brukes følgende tegn: > (større enn), < (mindre enn), == (er lik), != (ulik), >= (større enn eller lik), <= (mindre enn eller lik), | (eller), & (og).

Merk at en bruker kun enkel = med mindre det er i sammenheng med en if-betingelse.

Eksempel: `generate mann = 1 if kjønn == '1'`

Funksjonen `int()` gjør om tall til heltall (integer).

Funksjonen `substr()` trekker ut delstrenger fra variabler med verdier på tekstformat (alfanumerisk), der en angir hhv. variabelnavn, posisjon for første karakter, og antallet posisjoner som skal avleses. Eksempel på uttrekking av første karakter fra en tekstverdi:

```
pos1 = substr(varx,1,1)
```

Missing angis ved `sysmiss()`, der en spesifiserer variabelen inne i parentesene, f.eks. "drop if `sysmiss(inntekt)`" angir at enheter med manglende verdi på variabelen `inntekt` skal slettes.

Numeriske variabler: Verdier angis uten fnutter, f.eks. 1, 3, 10000

Alfanumeriske variabler: Verdier har tekstformat og angis med enkeltfnutter, f.eks. '1', '3', '10000' (dobbeltnutter kan også brukes).

4) Importere analysevariabler (fyll datasettet med variabler)

Du har nå en ferdig definert populasjon. Neste trinn blir å legge til variabler som du trenger for din analyse. Da bruker du kommandoen `import`.

De fleste data i databasen er organisert som forløpsdata, dvs. at variabelverdier oppdateres fortløpende. Valgfrie måletidspunkter angis da når en bruker `import`-kommandoen.

Faste opplysninger som kjønn, fødeland, fødselsdato etc: Måletidspunkter skal ikke oppgis

Enkelte variabler i databasen, såkalte tverrsnittsvariabler, er hentet fra statistiske kilder og måles kun ved faste datoer per år. Da må en oppgi en spesifikk måledato som blir foreslått av systemet når en prøver å taste inn dato i kommandolinjen. Du finner gyldige måledatoer i systemets variabeloversikter.

En fjerde variant er årlige akkumulerte verdier som viser den samme årsverdien uansett hvilken dato en velger innen et år. Dette gjelder hovedsakelig økonomiske opplysninger som inntekt, formue, gjeld etc. Det er da det samme hvilken dato en velger innenfor et aktuelt år.

Importer følgende variabler ved hjelp av kommandoen `import` (husk også å lage alias/kortnavn til hver enkelt variabel, f.eks. `import ds/BEFOLKNING_KJOENN as kjønn`):

- a) BEFOLKNING_KJOENN (kjønn, fast opplysning)
- b) BEFOLKNING_FODELAND (fødeland, fast opplysning)
- c) NUDB_BU per 2020-08-31 (6-sifret NUS-utdanningskode for høyeste fullførte utdanning, forløpsopplysning)
- d) BEFOLKNING_KOMMNR_FAKTISK per 2020-01-01 (4-sifret kommunekode for bosted, tverrsnittsopplysning)
- e) INNTEKT_WLONN per 2015-12-31 (årlig lønnsinntekt, akkumulert verdi => oppgitt verdi blir den samme uansett hvilken dato en velger innen det aktuelle året)
- f) INNTEKT_WLONN per 2016-12-31
- g) INNTEKT_WLONN per 2017-12-31
- h) INNTEKT_WLONN per 2018-12-31
- i) INNTEKT_WLONN per 2019-12-31

Legg merke til at kortnavnet til databasen legges til automatisk som prefiks ("ds/") på databasevariablene når du bruker import-kommandoen i det eksplorative kommandovinduet. Om du jobber i skripteditor må du huske å legge dette til selv.

NB! Personer med 0 i lønnsinntekt => missingverdi i databasen. Men dette er greit i denne analysen, siden vi kun er interesserte i å lage statistikk for positive inntekter. Om det er ønskelig å inkludere personer med null i inntekt, gjøres dette på følgende måte:

```
replace lønn = 0 if sysmiss(lønn)
```

5) Kjøre endimensjonal deskriptiv statistikk

Nå når vi har populasjonen og variablene klare, kan vi sette i gang og lage deskriptiv statistikk. En har en rekke hjelpemidler til rådighet.

For å hente ut statistiske tall for numeriske variabler bruker en vanligvis kommandoen `summarize`. For å hente ut frekvenser/opptellinger fordelt på ulike kategorier brukes kommandoen `tabulate`. En kan også kombinere `tabulate` og `summarize` for å vise statistiske tall fordelt på kategorier (f.eks. `tabulate kjønn, summarize(lønn)`).

Det er også mulig å lage grafiske fremvisninger som kakediagram (`piechart`), søylediagram (`barchart`), histogram (`histogram`), anonymisert plotdiagram (`hexbin`) og flytdiagram (`sankey`).

Deskriptiv statistikk kan enkelt kopieres og eksporteres inn i Excel- eller Google-regneark for videre bearbeiding og lagning av egne grafiske fremstillinger. På denne måten får en dessuten vist alle desimaler, og ikke bare standardfremvisningen i systemet. Grafisk output fra `microdata.no` kan kun eksporteres som bildefil og kan bare redigeres i bilderedigeringsprogram som "Paint".

- a) Vi starter med å lage en enkel tabell som viser prosentvis andel personer fordelt på kjønn. Kommandoen `tabulate` brukes til å lage tabeller med fordelinger av individer, enten antall eller prosentvis.

Forslag:

```
tabulate kjønn, cellpct
```

- b) La oss så lage en statistikk over lønnsutviklingen for vår populasjon. Kommandoen `summarize` brukes til å lage statistiske tall for numeriske variabler som f.eks. lønn. For å vise hvordan lønnsinntekten har endret seg i perioden 2015-2019, kan du liste opp alle variablene du vil måle i en enkelt `summarize`-kommando.

Forslag:

```
summarize lønn15 lønn16 lønn17 lønn18 lønn19
```

- c) Kommandoen `summarize` kan også brukes til å lage gini-koeffisienter som måler inntektsforskjeller i en gitt populasjon. Bruk opsjonen `gini` til dette. Opsjoner angis ved å legge til et komma bak kommando-uttrykket, etterfulgt av navnet på den aktuelle opsjonen.

Forslag:

```
summarize lønn15 lønn16 lønn17 lønn18 lønn19, gini
```

- d) Tallene kan fremstilles grafisk som søylediagram gjennom kommandoen `barchart`. Syntaxen har samme logikk ved at en lister opp variablene på samme måte som for `summarize`. I tillegg må en spesifisere hva slags statistikk en vil vise for søylediagrammene.

Forslag:

```
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19
```

I tillegg til gjennomsnitt (mean) kan du også vise antall individer med gyldige verdier (i dette tilfellet antall individer med positive verdier) gjennom å oppgi `count` i stedet for `mean`.

Forslag:

```
barchart (count) lønn15 lønn16 lønn17 lønn18 lønn19
```

Ved å angi `median` får du frem medianverdier i stedet. Prøv også dette!

- e) Studér også lønnsfordelingen i utvalget ditt for 2019. Bruk kommandoen `histogram` og benytt variabelen `lønn19`. Tips: Du kan bruke opsjonen `freq` for å vise antall som måleenhet på y-aksen i stedet for standardvisningen. Opsjonen `normal` kan dessuten brukes til å legge over en standard normalfordeling til referansebruk. Alle opsjoner kan brukes i kombinasjon med hverandre.

Forslag:

```
histogram lønn19, freq normal
```

- f) Lag en tilpasset lønnsfordeling der du i stedet for standardvisningen viser en fordeling på 10 søyler. Bruk opsjonen `bin(10)` til dette.

Lag også en lønnsfordeling der du bruker en inndeling i 4 søyler.

6) Kjøre todimensjonal deskriptiv statistikk

Det er mulig å lage statistikk og grafisk fremstilling med flere dimensjoner. Da trenger vi først å tilrettelegge variabler for dette formålet. Vi ønsker å se på hvordan lønn fordeler seg mellom kjønn, fødeland og utdanningsnivå. Kjønn er det ikke nødvendig å kode om, siden den kun har to kategorier. Men fødeland og utdanningsnivå derimot, trenger vi å lage grove kategorier for. Hvis ikke vil statistikken bli uoversiktlig.

Variabler med egne kategoriske inndelinger lages ved å bruke kommandoene `generate` og `replace`. Dette prinsippet gjelder både for dummyvariabler (med verdier 0 og 1) og andre variabler med flere enn to inndelinger. `generate` brukes til å sette en grunnverdi, og `replace` brukes til å modifisere verdien basert på spesifiserte if-kriterier. Se boksen under for eksempler.

Ved generering av kategoriske variabler, f.eks. dummyvariabler, er den enkleste metoden å først tilegne alle individer én av verdiene (bruk `generate`), for deretter å bruke `replace` til å endre verdi basert på gitte kriterier slik at klassifiseringen blir riktig. En kan bruke så mange `replace`-kommandoer en ønsker for å komme i mål med hele omkodingen av en gitt variabel, men bare én `generate`-kommando.

Et alternativ til `generate/replace` er å bruke den mer avanserte kommandoen `recode`. Merk at sistnevnte krever at variabler er på numerisk format.

Eksempel A:

```
generate mann = 1
replace mann = 0 if kjønn == '2'
```

Eksempel B:

```
generate norsk = 0
replace norsk = 1 if land == '000'
```

Eksempel C:

```
destring kjønn
recode kjønn (2 = 0)
rename kjønn mann
```

a)

La oss starte med å lage en dummyvariabel for fødeland: Norsk (=1) og ikke-norsk (=0). Kall variabelen "norsk". Pass på å alltid kode dummyvariabler som numeriske verdier (uten fnutter), ellers vil de ikke kunne benyttes i regresjonsanalyser.

Tips:

- i) Norsk \Leftrightarrow fødeland = '000'
- ii) Utenlandsk fødeland \Leftrightarrow øvrige koder inkludert missing (personer med midlertidig opphold)

Bruk kommandoen `tabulate` til å sjekke kodingen av dummyvariabelen "norsk". Bruk også opsjonen `cellpct` for prosentvis visning.

Prøv også kommandoen `piechart` for å lage et kakediagram over fordelingen norsk vs ikke-norsk.

b)

En enkelt måte å grovinndele utdanningsnivå på, er å trekke ut første siffer i den standardiserte NUS-koden (hierarkisk kode som går fra 0 (lavest) til 8 (høyest)).

Lag grovinndeling av utdanningsnivå (1. siffer-nivå):

Trekk ut første siffer i den 6-sifrede utdanningskoden ved å bruke funksjonen `substr()`.

Tips: `substr(utd,1,1)` trekker ut første siffer fra den alfanumeriske variabelen `utd`. Inputparameter nr. 2 angir startposisjon, mens parameter nr. 3 angir hvor mange posisjoner som skal avleses.

En kan lage alternative inndelinger av utdanningsnivå, i stedet for å bruke standard 1. siffer NUS-inndeling (= groveste NUS-nivå). Om en f.eks. ønsker 3 kategorier, kan en gjøre det slik:

```
import NUDB_BU 2011-06-01 as utd
generate utdnivå = substr(utd,1,1)
destring utdnivå, force
replace utdnivå = 1 if utdnivå <= 3
replace utdnivå = 2 if utdnivå == 4 | utdnivå == 5
replace utdnivå = 3 if utdnivå > 5
```

Merk at `NUDB_BU` er alfanumerisk. For å gjøre det lettere å omkode, brukes `destring-`kommandoen for å konverte verdiene til numerisk format (`force`-opsjonen brukes for å sikre at eventuelle verdier som inneholder andre tegn enn tall får verdien `sysmiss`). Da kan en bruke `<`, `>`, `<=`, `>=` etc, noe som ikke er mulig for alfanumeriske verdier. Husk at ledende 0-er blir fjernet ved numerisk konvertering. F.eks. '0301' blir til 301 ved konvertering til numerisk format.

Alternativt kan en bruke kommandoen `recode` i stedet for `generate/replace`:

```
import NUDB_BU 2011-06-01 as utd
generate utdnivå = substr(utd,1,1)
destring utdnivå, force
recode utdnivå (0/3 = 1) (4/5 = 2) (6/8 = 3)
```

Den siste linjen kan alternativt uttrykkes slik:

```
recode utdnivå (0 1 2 3 = 1) (4 5 = 2) (6 7 8 = 3)
```

Bruk kommandoen `tabulate` for å vise kodingen av utdanningsnivå. Bruk gjerne opsjonen `cellpct`.

Nå når vi har variablene klare, kan vi sette i gang med å lage flerdimensjonal statistikk.

For visning av frekvenser/antall fordelt på ulike kategorier, bruker en `tabulate` der en bare legger til flere enn én variabel i uttrykket. Jo flere variabler, jo flere dimensjoner. To variabler gir en toveis-tabell etc. Dette kan igjen kombineres med bruk av `summarize` som opsjon for visning av statistiske målinger i stedet for antall, f.eks. `tabulate kjønn sivilstand, summarize(lønn)`.

c)

Bruk kommandoen `tabulate` og lag krysstabeller som viser antall personer fordelt på variabler du har laget:

- "norsk" vs. kjønn
- utdanningsnivå vs. kjønn
- utdanningsnivå vs. "norsk"

For å lage krysstabeller, lister du bare opp mer enn én variabel i `tabulate`-uttrykket. Bruk dessuten opsjonen `rowpct` for å vise rekkeprosent i stedet for antall personer.

d)

Vi skal nå lage statistikk der vi ser nærmere på nivåer av utdanning. For å kunne rangere utdanningsnivåer (eller bruke dem i matematiske uttrykk), må vi først gjøre variabelen med den grove utdanningsinndelingen om til numerisk format. Bruk kommandoen `destring` til dette formålet, etterfulgt av navnet på variabelen du ønsker å gjøre numerisk.

Lag deretter oppsummerende statistikk over lønnsinntekt for årene 2015-2019, men bare for dem med utdanningsnivå < 2. Bruk kommandoen `summarize`.

Gjør det samme, men denne gang for dem med utdanningsnivå > 6.

e)

I stedet for å lage statistikk for undergrupper ved hjelp av if-betingelser, kan vi heller lage tabeller som viser statistikk fordelt på kategoriske variabler. Dette gir en mer oversiktlig visning av statistiske fordelinger. Dette gjøres ved å kombinere kommandoene `tabulate` og `summarize`. Ved å bruke `summarize` som en opsjon til `tabulate`, vises statistikk i stedet for antall personer i en tabell/krysstabell.

Vi skal først se på gjennomsnittlig lønnsinntekt for 2019 fordelt på kjønn.

Forslag:

```
tabulate kjønn, summarize(lønn19)
```

Lag tilsvarende tabeller, men fordel i stedet på hhv. fødeland (norsk) og utdanningsnivå (1-sifret). Bruk variablene du har laget.

Lag også en krysstabell der du ser på gjennomsnittlig lønnsinntekt for 2019 fordelt på utdanningsnivå OG kjønn.

f)

Lag et søylediagram som viser gjennomsnittlig lønnsinntekt over årene 2015-2019, fordelt på kjønn. Bruk kommandoen `barchart` med opsjonen `over()`.

Forslag:

```
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19, over(kjønn)
```

Gjør det samme, men fordel på hhv. fødeland (norsk) og utdanningsnivå.

g)

Også histogrammer kan fordeles på kategoriske variabler. Bruk kommandoen `histogram` og opsjonen `by()` til å vise lønnsfordelinger fordelt på hhv. kjønn og "norsk" i respektive fremvisninger.

Eksempel:

```
histogram lønn19, by(kjønn)
```

7) Ta vare på arbeidet du har gjort ved å lagre det som et skript

Et skript kan brukes til å foreta diverse justeringer i analysen. Alle kommandolinjene kan deretter kjøres på nytt i en bolck.

- Klikk på skript-knappen (knapp nr. 2 oppe i høyre hjørne) for å gå inn i skript-området
- Klikk på menyknappen helt opp til venstre
- Bla deg ned til og klikke på "Nytt program med historikk fra kommandolinjen"

The screenshot shows the RAIRO web interface. On the left, there is a sidebar with a menu. The main area is divided into a script editor on the left and a data table on the right. The script editor contains the following code:

```
»create-dataset drairdtest_diff
Et tomt datasett, drairdtest_diff, er opprettet og valgt

»import BOSTATTEFOT_BOSTED 2010-01-01 as dbokom100101, values( '1201' )
Importerte til drairdtest_diff med 256645 verdier

»import BEFOLKNING_REGSTAT 2010-01-01 as dregstat100101
Importerte til drairdtest_diff med 256645 verdier

»import BEFOLKNING_FOEDELSDATO as dfdato
Importerte til drairdtest_diff med 256645 verdier

»import BEFOLKNING_INVKAT as dinvkat
Importerte til drairdtest_diff med 256645 verdier

»import BEFOLKNING_KJODEN as dkjenn
Importerte til drairdtest_diff med 256645 verdier

»import SIVSTANDFDT_SIVSTAND 2010-01-01 as dsivst100101
Importerte til drairdtest_diff med 256645 verdier

»tabulate dinvkat dkjenn dsivst100101 if dinvkat == 'A', missing
```

The data table shows the following results:

	Mann	Kvinne	Sum
0 ukjent	6	--	6
Ugift	58531	52493	111036
Gift	35959	35257	71220
Enke/enkemann	2162	8785	10945
Skilt	6732	9499	16222
Separert	1214	1232	2447
Registrert partner	58	52	107

- d) Alt arbeidet ditt er nå lagt inn i et skript.
- e) Legg inn et navn på skriptet ditt i feltet like over skriptområdet. Skriptet ditt vil nå være lagret med dette navnet og kan finnes igjen ved å klikke på menyknappen oppe til venstre.
- f) Alternativt til a)-e) kan du bare skrive `save '.....'` i kommandolinjen, der du fyller inn et passende navn på skriptet inni enkelt-fnuttene. Deretter går du inn i skriptvinduet og jobber videre med skriptet ditt.
- g) Du kan nå kjøre arbeidet ditt samlet ved å klikke på “Kjør” eller “Send til kommandolinjen”, eventuelt gjøre justeringer først og deretter kjøre på nytt. Så lenge skriptet ditt har et definert navn, vil det bli autolagret når du jobber i det, på samme måte som i Google Docs.

Merk: Systemet husker tidligere kjørte kommandosekvenser. Så om du kjører samme skript to ganger på rad, så vil andre kjøring bare hente frem forrige resultat. Om du justerer på noe eller legger nye linjer til på slutten av skriptet, hentes de kjørte sekvenser frem fra minnet, og bare det nye kjøres. Endres derimot linjer helt i starten av et skript, anses alle sekvensene som nye, og alt kjøres på nytt igjen.

NB! Når et skript kjøres ved å bruke knappen “Send til kommandolinjen”, sendes alt til arbeidsområdet, og det som måtte ligge der fra før vil bli overskrevet av den nye kjøringen. Pass derfor på at arbeidet ditt i arbeidsområdet alltid blir tatt vare på via skript.

NB! En kan sikkerhetskopiere ved å kopiere skript over i egne dokumenter. Om en har tenkt å lime innholdet inn i et nytt skript etterpå, må en passe på å benytte et rent tekstformat via programmer som “Notisblokk” o.l. (tekst som kopieres inn i skripteditoren direkte fra programmer som Word o.l. kan skape krøll siden teksten da inneholder diverse formateringer som kan skape problemer ved kjøring av skript).

Er du i mål med oppgave 7? Da kan du prøve de litt mer avanserte oppgavene 8-10.

Du kan velge selv om du vil jobbe videre i skripteditoren eller i kommandolinjen i de neste oppgaver. Det anbefales å jobbe i skripteditoren og kjøre kommandoene derfra pga. muligheten for å lagre og systematisere arbeidet ditt. Editoren kan blant annet brukes til å kopiere og lime slik at en slipper å benytte så mye manuell programmering. Kommandolinjen derimot passer best for eksplorativt arbeid og for å gjøre seg kjent med variabler og syntax.

8) Lage finraffinert statistikk: Bruke labler som forklaring til koder i tabeller

Nå som vi behersker deskriptiv statistikk, kan vi gå et steg videre og se på hvordan vi kan gjøre tabellfremvisningene finere. Variabler en importerer fra databasen vil allerede ha verdilabler tilknyttet, og som viser hva de ulike kodene betyr. Men dette gjelder ikke variabler en lager selv ved bruk av `generate`-kommandoen.

a)

Vi skal lage verdilabler som vi knytter til variabelen vi lagde for å vise utdanningsnivå i oppgave 6.

Start med å bruke kommandoen `define-labels` etterfulgt av navnet på verdilabel-settet, og deretter alle verdiene og de tilhørende lablene.

Bruk disse lablene:

0	'Ingen utdanning'
1	Barneskole
2	Ungdomsskole
3	Videregående
4	'Videregående - avsluttende'
5	'Påbygging til videregående'
6	'UH-utdanning - lavere nivå'
7	'UH-utdanning - høyere nivå'
8	Forskerutdanning
9	Uoppgitt

Forslag:

```
define-labels utdlabel 0 'Ingen utdanning' 1 Barneskole 2 Ungdomsskole  
3 Videregående 4 'Videregående - avsluttende' 5 'Påbygging til  
videregående' 6 'UH-utdanning - lavere nivå' 7 'UH-utdanning - høyere  
nivå' 8 Forskerutdanning 9 Uoppgitt
```

Merk at labler som inneholder spesialtegn (mellomrom, bindestrek og alle andre tegn som ikke er bokstaver) må ha fnutter rundt. Du kan velge mellom enkeltfnutter og vanlige dobbeltfnutter. Verdiene må spesifiseres med det formatet de har i datasettet (i vårt tilfelle numeriske verdier som ikke skal ha fnutter rundt).

b)

Bruk kommandoen `assign-labels` til å knytte labelsettet du lagde med `utdanningsnivå`-variabelen. Kommandoen etterfølges av variabelnavnet og deretter navnet på labelsettet.

Forslag:

```
assign-labels utdnivå utdlabel
```

c)

Lag til slutt en krysstabell som viser antall personer fordelt på utdanningsnivå og kjønn. Bruk kommandoen `tabulate`. Kodene for utdanningsnivå skal nå vises med forklarings-teksten som du lagde i trinnene over. `Kjønn` har allerede kodelabler tilknyttet siden denne variabelen er hentet direkte fra databasen.

9) Lage tabeller med skreddersydde inndelinger

a)

Bruk kommandoene `generate` og `replace` til å lage en ny variabel som består av en enda grovere inndeling av utdanningsnivå. Kall variabelen f.eks. `utdgr`. Du gjør det på samme måte som i oppgave 6 a), bare at du bruker flere linjer med `replace` for å kode flere kategorier.

Forslag til grovinndeling av utdanningsnivå:

<u>Utdanningsnivå</u>	<u>Grovinndeling</u>	<u>Labeltekst</u>
0-5	1	Lav
6	2	Middels
7-8	3	Høy
9	9	Uoppgitt

Tips:

- Begynn med kommandoen `generate` og sett alle verdier for den nye variabelen `utdgr` til 1 dersom utdanningsnivå ≥ 0
- I neste linje bruker du `replace` til å erstatte med verdien 2 dersom utdanningsnivå ≥ 6
- Gjenta prosessen med en ny `replace`-kommando der du erstatter med verdien 3 dersom utdanningsnivå ≥ 7
- Til slutt erstatter du med verdien 9 dersom utdanningsnivå = 9 (uoppgitt)

Når du har gjort disse trinnene, vil du ha en ny variabel `utdgr` som består av tre utdanningsnivåer samt kategorien uoppgitt (9).

(Om dette blir for vanskelig, er det lov å titte litt i fasiten bakerst i oppgavesettet)

b)

Bruk kommandoen `define-labels` for å lage et nytt labelsett for den nye variabelen med grov inndeling, og deretter `assign-labels` for å koble lablene til variabelen `utdgr`, slik som i oppgave 8.

Lag til slutt en krysstabell ved hjelp av kommandoen `tabulate` som viser fordelingen mellom grovinndelt utdanningsnivå og kjønn. Bruk opsjonene `rowpct` og `freq` til å vise både rekkeprosent og antall i de ulike tabellcellene.

c)

Tabeller (og andre typer statistikk) vil som standard holde personer med missing-verdier utenfor når resultatet vises. Det kan i noen tilfeller være interessant å ta med også disse. Dette gjøres ved å bruke opsjonen `missing`.

Lag en tilsvarende tabell som i b) der du inkluderer manglende verdier. Bruk opsjonen `missing`, gjerne i kombinasjon med `rowpct` og `freq`. Du vil nå se at det finnes en liten gruppe individer som ikke har verdi på høyeste fullførte utdanning (dette er vanligvis personer med utenlandsk bakgrunn).

10) Kjøre en enkel regresjonsanalyse

Gjennomfør en lineær regresjonsanalyse der en analyserer effekter av ulike bakgrunns karakteristika på lønnsinntekt i 2019 ved å bruke kommandoen `regress`. Benytt variabler du har laget tidligere i oppgavesettet:

- Lønnsinntekt per 2019 (avhengig variabel)
- Alder
- Norsk
- Mann (lag en dummyvariabel for mann ut i fra variabelen `kjønn`)
- Oslo (lag en dummyvariabel som settes lik 1 dersom `bosted == '0301'`)
- Høy utdanning (lag en dummyvariabel for høyt utdanningsnivå ut i fra variabel for høyeste fullførte utdanning, nivå $\geq 7 \Rightarrow$ verdien 1, resten 0)

NB! Husk at den avhengige variabelen skal listes først i `regress`-uttrykket.

Forslag:

```
generate mann = 0
replace mann = 1 if kjønn == '1'

generate oslo = 1 if bosted == '0301'
replace oslo = 0 if bosted != '0301'
tabulate oslo, cellpct

generate høy_utd = 1 if utdnivå >= 7
replace høy_utd = 0 if utdnivå >= 0 & utdnivå < 7

regress lønn19 alder norsk oslo mann høy_utd
```

Fasit (syntax):

1)

```
require no.ssb.fdb:14 as ds
```

2)

```
create-dataset totalpop
```

3a)

```
import ds/BEFOLKNING_FOEDSELS_AAR_MND as faarmnd  
generate alder = 2020 - int(faarmnd/100)
```

3b)

```
histogram alder, discrete
```

3c)

```
import ds/BEFOLKNING_STATUSKODE 2020-01-01 as regstat  
keep if regstat == '1'  
histogram alder, discrete  
summarize alder
```

3d)

```
keep if alder > 30 & alder < 50
```

4)

```
import ds/BEFOLKNING_KJOENN as kjønn  
import ds/BEFOLKNING_FODELAND as land  
import ds/NUDB_BU 2020-08-31 as utd  
import ds/BEFOLKNING_KOMMNR_FAKTISK 2020-01-01 as bosted  
import ds/INNTEKT_WLONN 2015-12-31 as lønn15  
import ds/INNTEKT_WLONN 2016-12-31 as lønn16  
import ds/INNTEKT_WLONN 2017-12-31 as lønn17  
import ds/INNTEKT_WLONN 2018-12-31 as lønn18  
import ds/INNTEKT_WLONN 2019-12-31 as lønn19
```

5a)

```
tabulate kjønn, cellpct
```

5b)

```
summarize lønn15 lønn16 lønn17 lønn18 lønn19
```

5c)

```
summarize lønn15 lønn16 lønn17 lønn18 lønn19, gini
```

5d)

```
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19  
barchart (count) lønn15 lønn16 lønn17 lønn18 lønn19  
barchart (median) lønn15 lønn16 lønn17 lønn18 lønn19
```

5e)

```
histogram lønn19, freq  
histogram lønn19, freq normal
```

5f)

```
histogram lønn19, bin(10)  
histogram lønn19, bin(4)
```

6a)

```
generate norsk = 0  
replace norsk = 1 if land == '000'  
tabulate norsk  
tabulate norsk, cellpct  
piechart norsk
```

6b)

```
generate utdnivå = substr(utd,1,1)  
tabulate utdnivå, cellpct
```

6c)

```
tabulate norsk kjønn, rowpct  
tabulate utdnivå kjønn, rowpct  
tabulate utdnivå norsk, rowpct
```

6d)

```
destring utdnivå  
summarize lønn15 lønn16 lønn17 lønn18 lønn19 if utdnivå < 2  
summarize lønn15 lønn16 lønn17 lønn18 lønn19 if utdnivå > 6
```

6e)

```
tabulate kjønn, summarize(lønn19)  
tabulate norsk, summarize(lønn19)  
tabulate utdnivå, summarize(lønn19)  
tabulate utdnivå kjønn, summarize(lønn19)
```

6f)

```
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19, over(kjønn)
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19, over(norsk)
barchart (mean) lønn15 lønn16 lønn17 lønn18 lønn19, over(utdnivå)
```

6g)

```
histogram lønn19, by(kjønn)
histogram lønn19, by(norsk)
```

8a)

```
define-labels utdlabel 0 'Ingen utdanning' 1 Barneskole 2 Ungdomsskole 3
Videregående 4 'Videregående - avsluttende' 5 'Påbygging til videregående' 6
'UH-utdanning - lavere nivå' 7 'UH-utdanning - høyere nivå' 8
Forskerutdanning 9 Uoppgitt
```

8b)

```
assign-labels utdnivå utdlabel
```

8c)

```
tabulate utdnivå kjønn
```

9a)

```
generate utdgr = 1 if utdnivå >= 0
replace utdgr = 2 if utdnivå >= 6
replace utdgr = 3 if utdnivå >= 7
replace utdgr = 9 if utdnivå == 9
```

9b)

```
define-labels utdlabel2 1 Lav 2 Middels 3 Høy 9 Uoppgitt
assign-labels utdgr utdlabel2
tabulate utdgr kjønn, rowpct freq
```

9c)

```
tabulate utdgr kjønn, rowpct freq missing
```

10)

```
generate mann = 0
replace mann = 1 if kjønn == '1'
```

```
generate oslo = 1 if bosted == '0301'  
replace oslo = 0 if bosted != '0301'  
tabulate oslo, cellpct
```

```
generate høy_utd = 1 if utdnivå >= 7  
replace høy_utd = 0 if utdnivå >= 0 & utdnivå < 7
```

```
regress lønn19 alder norsk oslo mann høy_utd
```
