

Case 4

277

Instruksjoner

Denne oppgaven skal løses interaktivt i RStudio ved å legge inn egen kode og kommentarer. Det ferdige dokumentet lagres med kandidatnummeret som navn `[kandidatnummer]_SOK1004_C4_H22.qmd` og lastes opp på deres GitHub-side. Hvis du har kandidatnummer 43, så vil filen hete `43_SOK1004_C4_H22.qmd`. Påse at koden kjører og at dere kan eksportere besvarelsen til pdf. Lever så lenken til GitHub-repositoriet i Canvas.

Bakgrunn, læringsmål

Innovasjon er en kilde til økonomisk vekst. I denne oppgaven skal vi se undersøke hva som kjennetegner bedriftene som bruker ressurser på forskning og utvikling (FoU). Dere vil undersøke FoU-kostnader i bedriftene fordelt på næring, antall ansatte, og utgiftskategori. Gjennom arbeidet vil dere repetere på innhold fra tidligere oppgaver og øve på å presentere fordelinger av data med flere nivå av kategoriske egenskaper.

Last inn pakker

```
# output | false
rm(list=ls())
library(tidyverse)
library(rjstat)
library(gdata)
library(httr)
```

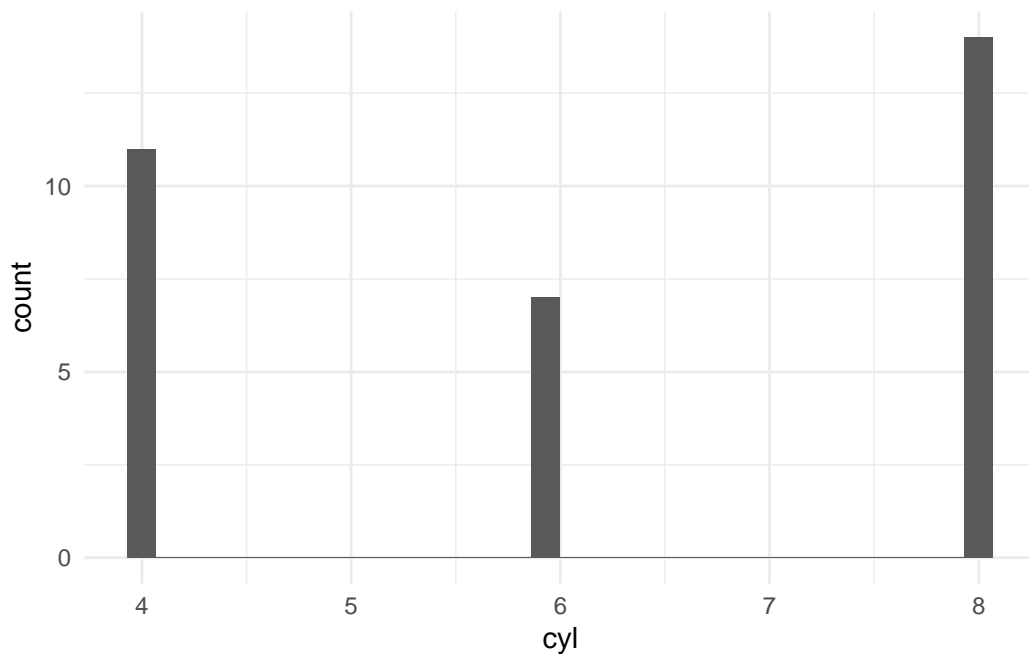
Oppgave I: Introduksjon til histogram

Et histogram eller frekvensfordeling er en figur som viser hvor ofte forskjellige verdier oppstår i et datasett. Frekvensfordelinger spiller en grunnleggende rolle i statistisk teori og modeller. Det er avgjørende å forstå de godt. En kort innføring følger.

La oss se på et eksempel. I datasettet `mtcars` viser variabelen `cyl` antall sylindere i motorene til kjøretøyene i utvalget.

```
data(mtcars)
mtcars %>%
  ggplot(aes(cyl)) +
  geom_histogram() +
  theme_minimal()
```

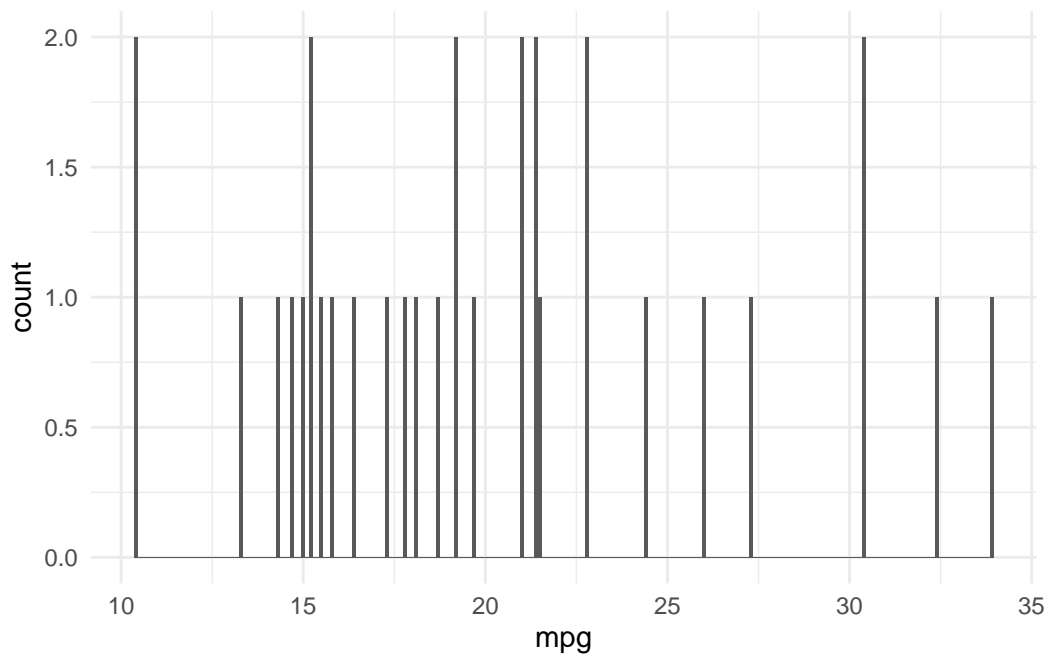
``stat_bin()`` using ``bins = 30``. Pick better value with ``binwidth``.



Verdiene av variabelen er gitt ved den horisontale aksene, antall observasjoner på den vertikale aksene. Vi ser at det er 11, 7, og 14 biler med henholdsvis 4, 6, og 8 sylindere.

La oss betrakte et eksempel til. Variabelen `mpg` i `mtcars` måler gjennomsnittlig drivstofforbruk i uanstendige engelske enheter. Variabelen er målt med ett desimal i presisjon.

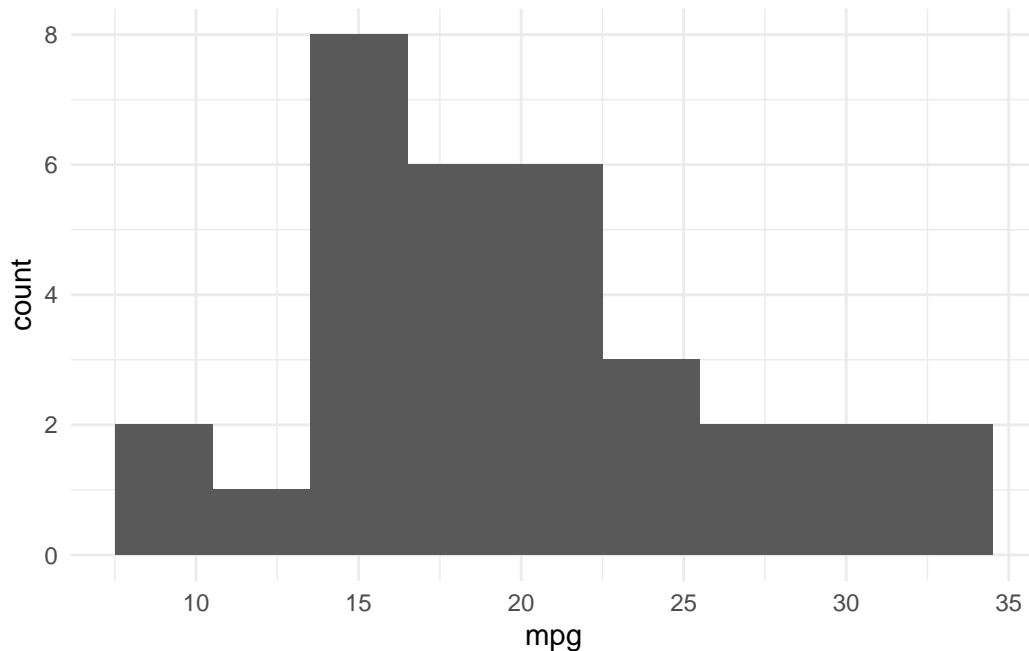
```
data(mtcars)
mtcars %>%
  ggplot(aes(mpg)) +
  geom_histogram(binwidth=0.1) +
  theme_minimal()
```



Datasettet inneholder mange unike verdier, hvilket gir utslag i et flatt histogram, noe som er lite informativt. Løsningen da er å gruppere verdier som ligger i nærheten av hverandre. Kommandoen `binwidth` i `geom_histogram()` bestemmer bredden av intervallene som blir slått sammen. Kan du forklare hvorfor alle unike verdier blir telt ved å bruke `binwidth = 0.1`?

Eksperimenter med forskjellige verdier for `binwidth` og forklar hva som kjennetegner en god verdi.

```
# løs oppgave I her
data(mtcars)
mtcars %>%
  ggplot(aes(mpg)) +
  geom_histogram(binwidth=3) +
  theme_minimal()
```



En god verdi på binwidth er en verdi der vi klarer å lese av representativ data fra den samlingen av data. Det vil si at dataen som befinner seg i den binwidth'en ikke har for stor spredning til at dataen vi kan lese av ikke blir representativ for dataen i binwidth'en.

Oppgave II: Last ned og rydd i data

Vi skal nå undersøke dataene i [Tabell 07967: Kostnader til egenutført FoU-aktivitet i næringslivet, etter næring \(SN2007\) og sysselsettingsgruppe \(mill. kr\) 2007 - 2020 SSB](#). Dere skal laste de ned ved hjelp av API. Se [brukerveiledningen](#) her.

Bruk en JSON-spørring til å laste ned alle statistikkvariable for alle år, næringer, og sysselsettingsgrupper med 10-19, 20-49, 50-99, 100-199, 200 - 499, og 500 eller flere ansatte. Lagre FoU-kostnader i milliarder kroner. Sørg for at alle variabler har riktig format, og gi de gjerne enklere navn og verdier der det passer.

Hint. Bruk lenken til SSB for å hente riktig JSON-spørring og tilpass koden fra case 3

```
# besvar oppgave II her

# Henter data fra tabellen til SSB.
url <- "https://data.ssb.no/api/v0/no/table/07967/"
```

```
# Klipp og lim spørringen fra SSB.
```

```
query <- '{  
  "query": [  
    {  
      "code": "NACE2007",  
      "selection": {  
        "filter": "item",  
        "values": [  
          "A-N",  
          "C",  
          "G-N",  
          "A-B_D-F"  
        ]  
      }  
    },  
    {  
      "code": "SyssGrp",  
      "selection": {  
        "filter": "item",  
        "values": [  
          "NyAlle",  
          "05-09",  
          "Alle",  
          "10-19",  
          "20-49",  
          "10-49",  
          "50-99",  
          "100-199",  
          "200-499",  
          "500+"  
        ]  
      }  
    }  
  ],  
  "response": {  
    "format": "json-stat2"  
  }  
}'
```

```
# Bruker samme framgangsmåte som i case 3.
```

```
hent_indeks.tmp <- url %>%
```

```

POST(body = query, encode = "json")

df <- hent_indeks.tmp %>%
  content("text") %>%
  fromJSONstat() %>%
  as_tibble()

# Gir variablene nye navn.
df <- df %>%
  rename("verdi" = "value") %>%
  rename("variabel" = "statistikkvariabel") %>%
  rename("gruppe" = "sysselsettingsgruppe") %>%
  rename("næring" = "næring (SN2007)")

# Deler verdiene på 1000 slik at de vises i milliarder(1000 mill. / 1000 = 1 mrd.).
df <- df %>%
  mutate(verdi = verdi/1000)

```

Oppgave III: Undersøk fordelingen

Vi begrenser analysen til bedrifter med minst 20 ansatte og tall fra 2015 - 2020. Lag en figur som illustrerer fordelingen av totale FoU-kostnader fordelt på type næring (industri, tjenesteyting, andre) og antall ansatte i bedriften (20-49, 50-99, 100-199, 200-499, 500 og over). Tidsdimensjonen er ikke vesentlig, så bruk gjerne histogram.

Merknad. Utfordringen med denne oppgaven er at fordelingene er betinget på verdien av to variable. Kommandoen `facet_grid()` kan være nyttig til å slå sammen flere figurer på en ryddig måte.

```

# besvar oppgave III her

# Filtrerer ut alle tall fra før 2015, filtrerer ut alle næringer og filtrerer ut alle var
df_copy_1 <- df %>%
  filter(år > 2014) %>%
  filter(variabel == "FoU-kostnader i alt") %>%
  filter(næring != "Alle næringer") %>%
  filter(gruppe != "10-19 sysselsatte") %>%
  filter(gruppe != "10-49 sysselsatte") %>%
  filter(gruppe != "5-9 sysselsatte") %>%
  filter(gruppe != "Alle (minst 10 sysselsatte)") %>%

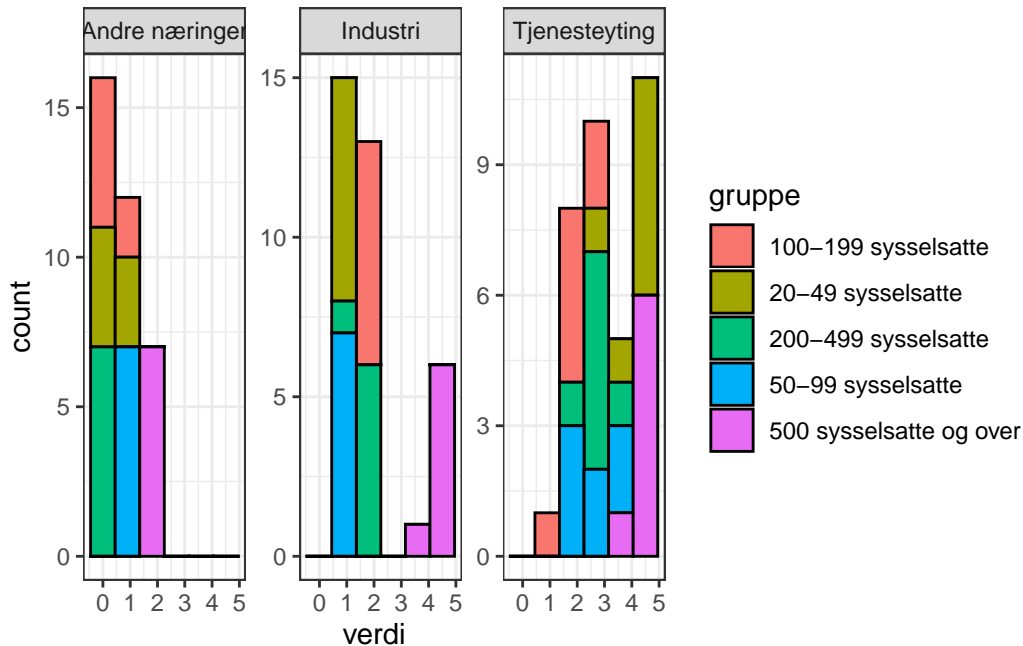
```

```

filter(gruppe != "Alle (minst 5 sysselsatte)")

# Plotter figurene med histogram.
df_copy_1 %>%
  ggplot(aes(x = verdi, fill = gruppe)) +
  geom_histogram(color = "black", binwidth = 0.9) +
  facet_wrap(~ næring, scale = "free_y") +
  theme_bw()

```



Oppgave IV: Undersøk fordelingen igjen

Kan du modifisere koden fra oppgave II til å i tillegg illustrere fordelingen av FoU-bruken på lønn, innleie av personale, investering, og andre kostnader?

Merknad. Kommandoen `fill = [statistikkvariabel]` kan brukes i et histogram.

```

# besvar oppgave III her

# Filtreer ut alle tall fra før 2015, filtrerer ut alle næringer og filtrerer ut alle var
df_copy_2 <- df %>%
  filter(år > 2014) %>%

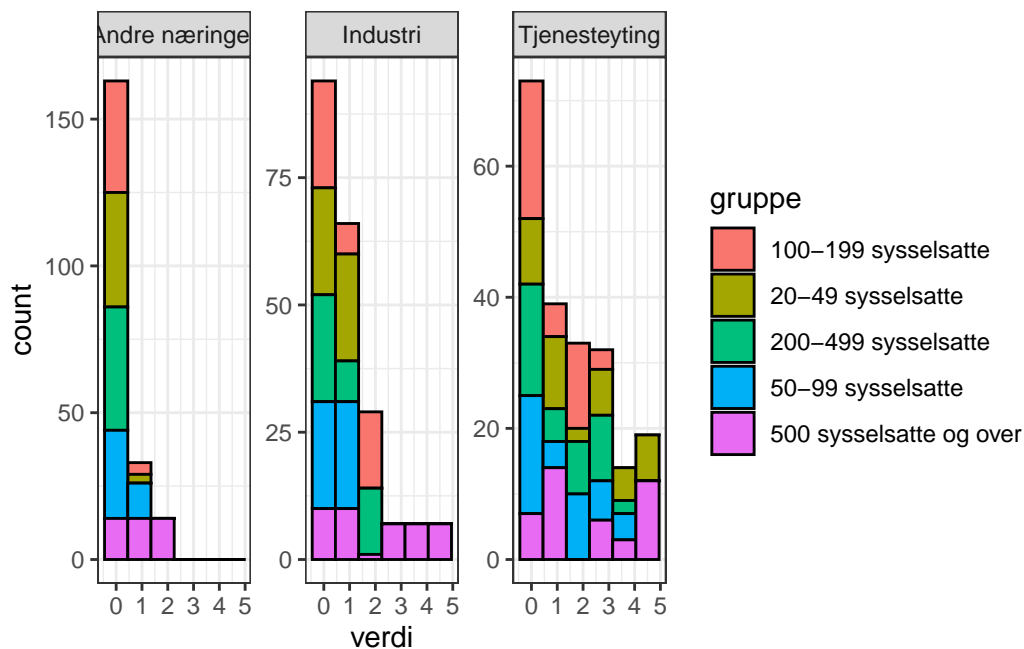
```

```

filter(næring != "Alle næringer") %>%
filter(gruppe != "10-19 sysselsatte") %>%
filter(gruppe != "10-49 sysselsatte") %>%
filter(gruppe != "5-9 sysselsatte") %>%
filter(gruppe != "Alle (minst 10 sysselsatte)") %>%
filter(gruppe != "Alle (minst 5 sysselsatte)")

# Plotter figurene med histogram.
df_copy_2 %>%
  ggplot(aes(x = verdi, fill = gruppe)) +
  geom_histogram(color = "black", binwidth = 0.9) +
  facet_wrap(~ næring, scale = "free_y") +
  theme_bw()

```



Kildeliste

- Daniel D. (06.05.2015). GGplot lessons 3.1: Facet Wrap [Video]. Youtube: <https://www.youtube.com/watch?v=wGDdJo1qPXM>