



på



GRUNNKURS

SUSIE JENTOFT

DAG 2



**Statistisk sentralbyrå**  
Statistics Norway

# Mål

- Kjennskap til programmet R og RStudio
- Åpne RStudio og kjøre enkle beregninger
- Lese inn data
- Behandle data
- Lage tabeller og oppsummere
- Lage figurer



# Agenda

## **Dag 2: 2. februar**

14:00 – Oppsummering og løsninger

14:10 – Databehandling og tabeller

14:30 – Øvelse 3

15:00 – Kobling og figurer

15:30 – Øvelse 4

15:50 – Oppsummering



# Øvelser gjennomgang



# Data behandling med tidyverse

- Gjør koden ryddigere
- Pipelines %>%

Base R:

```
leave_house(get_dressed(get_out_of_bed(wake_up(me))))
```

tidyverse:

```
me %>%  
  wake_up() %>%  
  get_out_of_bed() %>%  
  get_dressed() %>%  
  leave_house()
```

# Data behandling med tidyverse

Funksjon	Hva den gjøre
mutate()	Lage nye variable
rename()	Gi en variabel et nytt navn
filter()	Selektere noen enheter/rad
select()	Selektere noen variable
summarise()	Oppsummering av en variabel
group_by()	Gjøre prosesser innenfor grupper



# Lage nye variabler: mutate()

- Kan brukes som en del av en pipeline

```
datanavn %>%  
  mutate(nyvariabel = 1000)
```

Gir variabel et navn

```
datanavn %>%  
  mutate(nyvariabel = oldvariabel * 1000)
```

Gir variabel et navn

Eksisterende variabel

+ hva skal gjøres



**Statistisk sentralbyrå**  
Statistics Norway

# Endre variabelnavn: rename()

```
datanavn %>%  
  rename(nyttnavn = gammeltnavn)
```





# Velg noen rader: filter()

- For å velge ut noen rader bruker vi filter()
- Skriv logiske setning inn i parentes.
- Flere logiske setninger kan brukes sammen (skille med , )

```
datanavn %>%  
  filter(condition)
```

- Igjen: Ingenting lagres uten <-



# Velg ut noen variabler: select()

- Brukes med pipelines
- Skriv variabelnavn i parentes
- En eller flere variabler (skille med , )
- Brukes sammen med andre funksjoner (for eks. filter( ) )

```
datanavn %>%  
  select(variabelnavn)
```

```
datanavn %>%  
  filter(condition) %>%  
  select(variabelnavn)
```



# Oppsummering/aggregering: summarise()

- Ta oppsummering (summen, gjennomsnitt, median, antall) av en variabel med summarise()

```
datanavn %>%  
  summarise(oppsummeringsnavn = mean(variabelnavn))
```

Gir oppsummering et navn

Hva skal gjøres:

- mean()
- median()
- sum()
- n()

Eksisterende variabel



**Statistisk sentralbyrå**  
Statistics Norway

# Gruppering: group\_by()

- Gjøre alle prosesser etterpå innen hver gruppe

```
datanavn %>%  
  group_by(grupperingsvariabel) %>%  
  summarise(oppsummeringsnavn = mean(variabelnavn))
```



# Eksempler på dapla



# Øvelser 3

- Øvelsene til oppgavesett 3 er på fil: **øvelser\_dag2.R**



DataA

ID	Kommune	Omsetning
1	0301	212400
2	0301	75000
3	5011	620000

DataB

ID	Antall_ansatte
2	8
3	22
4	3

left\_join(DataA, DataB)

ID	Kommune	Omsetning	Antall ansatte
1	0301	212400	NA
2	0301	75000	8
3	5011	620000	22

right\_join(DataA, DataB)

ID	Kommune	Omsetning	Antall ansatte
2	0301	75000	8
3	5011	620000	22
4	NA	NA	3

# Koble to datasett

DataA

ID	Kommune	Omsetning
1	0301	212400
2	0301	75000
3	5011	620000

DataB

ID	Antall_ansatte
2	8
3	22
4	3

inner\_join(DataA, DataB)

ID	Kommune	Omsetning	Antall ansatte
2	0301	75000	8
3	5011	620000	22

full\_join(DataA, DataB)

ID	Kommune	Omsetning	Antall ansatte
1	0301	212400	NA
2	0301	75000	8
3	5011	620000	22
4	NA	NA	3

# Koble to datasett



# Koble to datasett

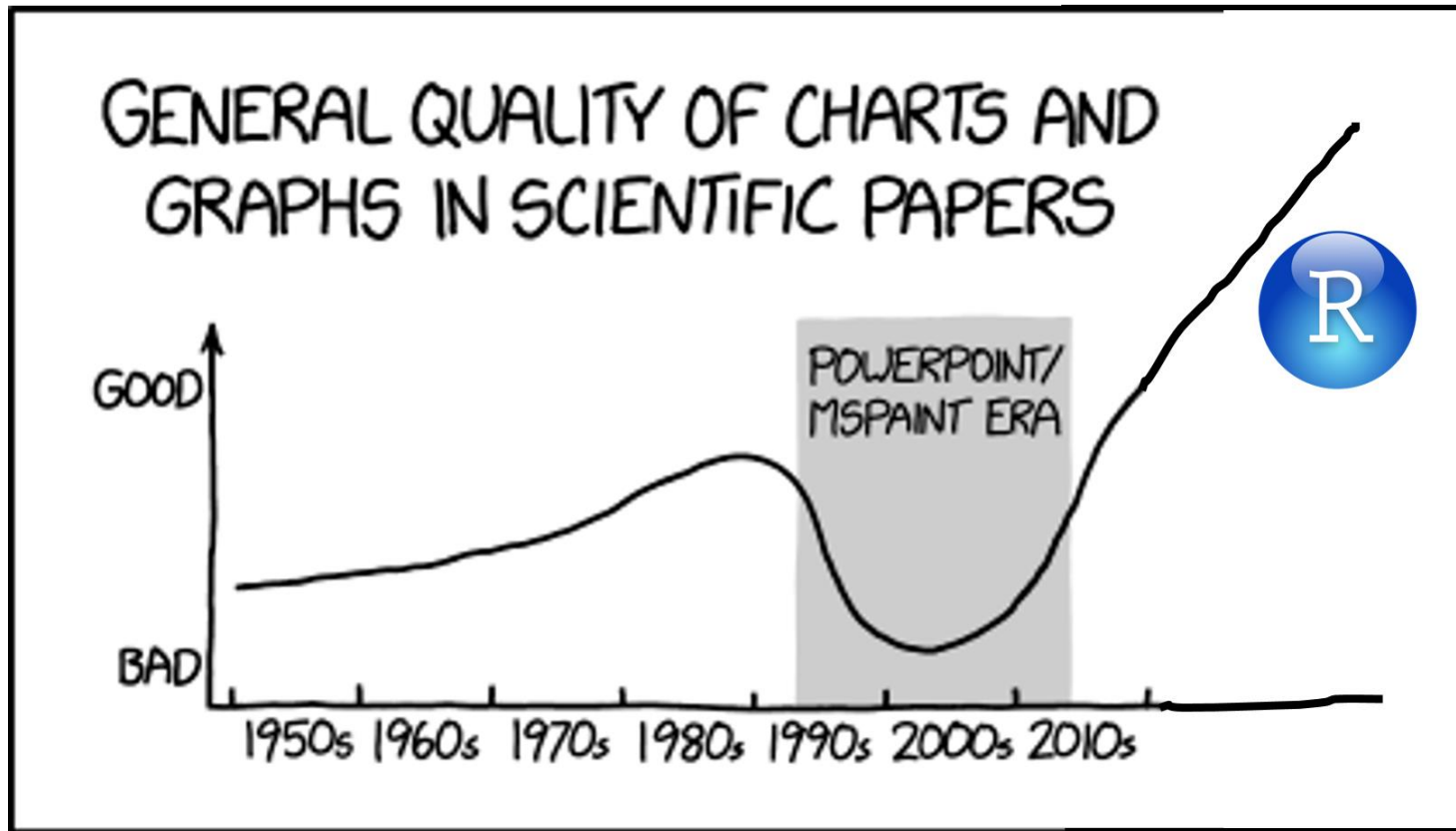
- Bruk **by** = for å spesifisere nøkkel variabel til å koble på

```
by = c("aar" = "year")
```

```
kobletdata <- left_join(datanavn1, datanavn2, by = variabelnavn)
```

- Flere variabler kan brukes for å koble på (som en vektor)

# Plotting



# Plotting med ggplot()

- **aes** : aesthetics, hvilke variabler
- **geom\_** : hva slags figur
- **stat** : hva slags statistisk aggregat å presentere

**Table 18-1 A Selection of Geoms and Associated Default Stats**

<i>Geom</i>	<i>Description</i>	<i>Default Stat</i>
<code>geom_bar()</code>	Bar chart	<code>stat_bin()</code>
<code>geom_point()</code>	Scatterplot	<code>stat_identity()</code>
<code>geom_line()</code>	Line diagram, connecting observations in order by x-value	<code>stat_identity()</code>
<code>geom_boxplot</code>	Box-and-whisker plot	<code>stat_boxplot()</code>
<code>geom_path</code>	Line diagram, connecting observations in original order	<code>stat_identity()</code>
<code>geom_smooth</code>	Add a smoothed conditioned mean	<code>stat_smooth()</code>
<code>geom_histogram</code>	An alias for <code>geom_bar()</code> and <code>stat_bin()</code>	<code>stat_bin()</code>



# Søylediagram

```
ggplot(aes(variabelnavn)) +  
  geom_bar()
```

Bruke + for å legge til figurtype

Spesifisere variabelen

Spesifisere søylediagram

```
ggplot(aes=c(x = variabelnavn1, y = variabelnavn2)) +  
  geom_bar(stat="identity")
```

Spesifisere x og y variablene

Spesifisere å bruke verdi



# Eksempler på dapla



# Øvelse 4

- Øvelsene til oppgavesett 1 er på fil: **øvelser\_4.R**



# Videre

- sparklyr: <https://spark.rstudio.com/>



# Oppsummering

- Husk library( )
- Les inn filer: read\_csv2( ) read\_sas( )
- Ny variabel: mutate( )
- Velg noen linje: filter( )
- Aggregere/oppsummere: summarise( )
- Figur: ggplot( ), aes( ), geom\_...( )
- <https://wiki.ssb.no/display/s880/For+R+brukere>
- Yammer: R i SSB
- Google

