

SQL-Datenbanken abfragen und Versionierung mit git

Unit 6

Ziele für heute

1. eine Verbindung zu einer Datenbank herstellen und SQL-Abfragen mit `{DBI}` ausführen
2. den grundlegenden Workflow von Git in RStudio erklären
3. die wichtigsten Konzepte des Kurses zusammenfassen

Praktikum 06a

CO2-Emissionen

30:00
rstatsBL - Data Science mit R



SQL Database

SQL Database

- **DBI + RPostgres** → Verbindung mit SQL Datenbank
- **dbplyr** → **dplyr**-Verben mit Datenbank

```
1 library(DBI) # dbConnect()  
2 library(RPostgres) # Postgres()
```

SQL Database

- **DBI + RPostgres** → Verbindung mit SQL Datenbank
- **dbplyr** → **dplyr**-Verben mit Datenbank

```
1 library(DBI) # dbConnect()
2 library(RPostgres) # Postgres()
3
4 con <- dbConnect(
5   Postgres(),
6   sslmode = "allow",
7   host = "host-name",
8   port = port-number,
9   dbname = "db-name",
10  user = "username",
11  password = "my_password"
12 )
```



Niemals Passwort im Skript speichern!

SQL Database

- DBI + RPostgres → Verbindung mit SQL Datenbank
- dbplyr → dplyr-Verben mit Datenbank

```
1 library(DBI) # dbConnect()
2 library(RPostgres) # Postgres()
3
4 con <- dbConnect(
5   Postgres(),
6   sslmode = "allow",
7   host = "host-name",
8   port = port-number,
9   dbname = "db-name",
10  user = "username",
11  password = rstudioapi::askForPassword("Datenbank Passwort")
12 )
```

config.yml

```

1 default:
2   DWH_var:
3     host: "host-name"
4     user: "username"
5     password: "DWH_pass"
6     port: port-number
7     dbname: "db-name"

```

```

1 library(DBI) # dbConnect()
2 library(RPostgres) # Postgres()
3
4 con <- dbConnect(
5   Postgres(),
6   sslmode = "allow",
7   host = config::get("DWH_var")$host,
8   dbname = config::get("DWH_var")$dbname,
9   user = config::get("DWH_var")$user,
10  password = config::get("DWH_var")$password,
11  port = config::get("DWH_var")$port
12 )

```

config.yml

```
1 default:  
2   DWH_var:  
3     host: "host-name1"  
4     user: "username"  
5     password: "DWH_pass"  
6     port: port-number1  
7     dbname: "db-name1"  
8   DWH_hrm:  
9     host: "host-name"  
10    user: "username"  
11    password: "DWH_pass"  
12    port: port-number2  
13    dbname: "db-name2"  
14  
15  andere_DB:  
16    DB_fin:  
17    host: "host-name3"  
18    user: "username"
```

SQL Daten Bearbeiten

```
1 con
```

```
1 dbListObjects(con)
```

SQL Daten Bearbeiten

```
1 dbListObjects(con, Id(schema = "schema_name"))
```

SQL Daten Bearbeiten

```
1 library(tidyverse)
2 df_db <- tbl(con, Id(schema = "schema_name", table = "table_name"))
3 class(df_db)
1 glimpse(df_db, width = 20)
```

SQL Daten Bearbeiten

```
1 df_db_small <- df_db |>
2   select(var1, district_k, var2) |>
3   mutate(
4     var1 = as.character(var1), # dbplyr cannot translate factor() to sql
5     district_k = as.character(district_k),
6     district_k = case_when(
7       district_k == "1301" ~ "Arlesheim",
8       district_k == "1302" ~ "Laufen",
9       district_k == "1303" ~ "Liestal",
10      district_k == "1304" ~ "Sissach",
11      .default = "Waldenburg"
12    )
13  ) |>
14  summarise(var3 = sum(var2), .by = c(var1, district_k))
```

SQL Daten Bearbeiten

```
1 df_db_small |>
2   ggplot(
3     aes(
4       x = var1,
5       y = var3,
6       fill = fct_reorder(
7         district_k,
8         var3
9       )
10      )
11    ) +
12   geom_col() +
13   labs(
14     x = "",
15     y = "Total",
16     fill = "Bezirk"
17   ) +
18   scale_v_continuous()
```

SQL Query: `show_query()`

```
1 df_db_small |>  
2   show_query()
```

SQL Queries in Quarto Dateien

```
1  ```{sql}
2 #| connection: con
3 #| output.var: "db_sql_small"
4 #| code-line-numbers: "1-4|"
5
6 SELECT "var1", "district_k", SUM("var2") AS "var3"
7 FROM (
8   SELECT
9     "var1",
10    CASE
11      WHEN ("district_k" = '1301') THEN 'Arlesheim'
12      WHEN ("district_k" = '1302') THEN 'Laufen'
13      WHEN ("district_k" = '1303') THEN 'Liestal'
14      WHEN ("district_k" = '1304') THEN 'Sissach'
15      ELSE 'Waldenburg'
16    END AS "district_k",
17    ""
18  FROM (
```

Mit Resultierenden Daten Weiter Arbeiten

▼ Code

```
1 db_sql_small |>
2   ggplot(
3     aes(
4       x = var1,
5       y = var3,
6       fill = fct_reorder(district_k, var3)
7     )
8   ) +
9   geom_col() +
10  labs(
11    x = "",
12    y = "Total",
13    fill = "Bezirk"
14  ) +
15  scale_y_continuous(
16    label = scales::label_number(prefix = "CHF ", big.mark = ""))
17 ) +
18 theme_minimal()
```

Resultierende Daten aus Datenbank Holen

```
1 class(df_db_small)
```

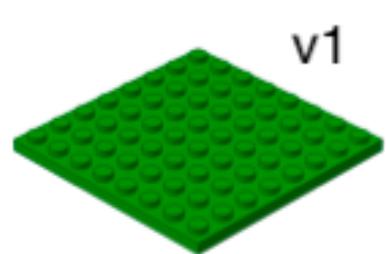
👉 Daten in DB (*lazy loading*)

```
1 df <- df_db_small |>
2   collect()
3
4 class(df)
```

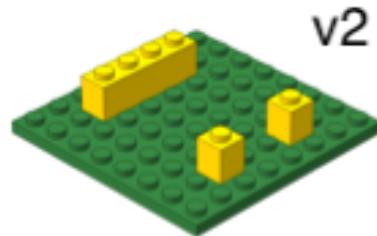
👉 Daten im Dataframe!

Versionierung und Kollaboration: git

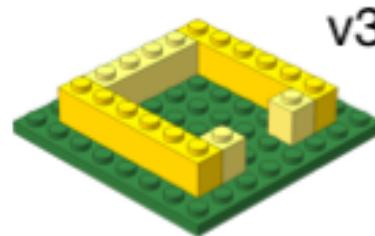
Versionierung



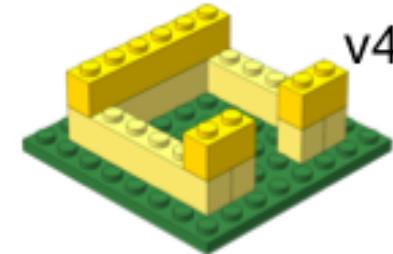
v1



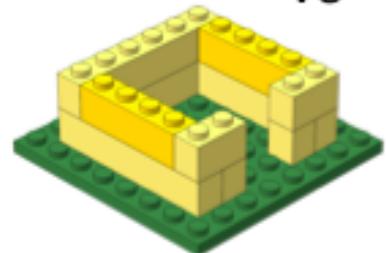
v2



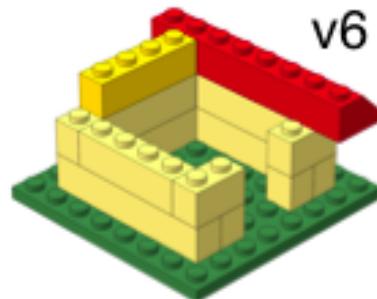
v3



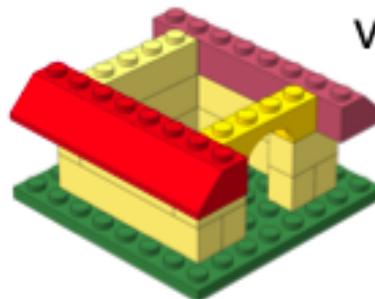
v4



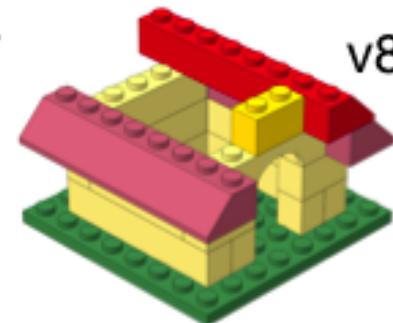
v5



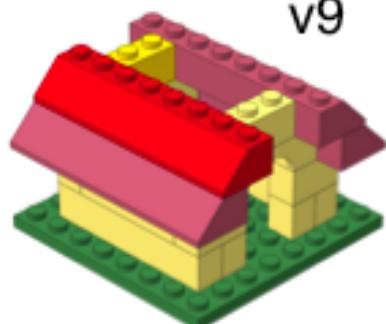
v6



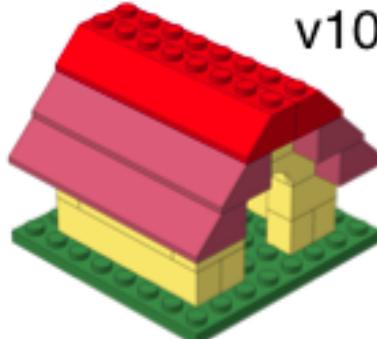
v7



v8



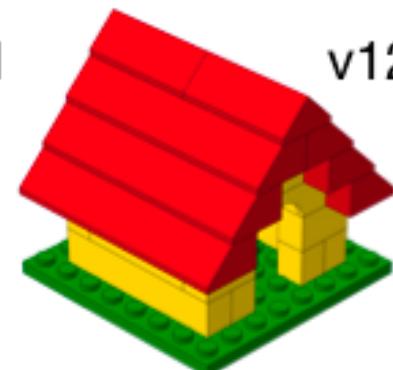
v9



v10

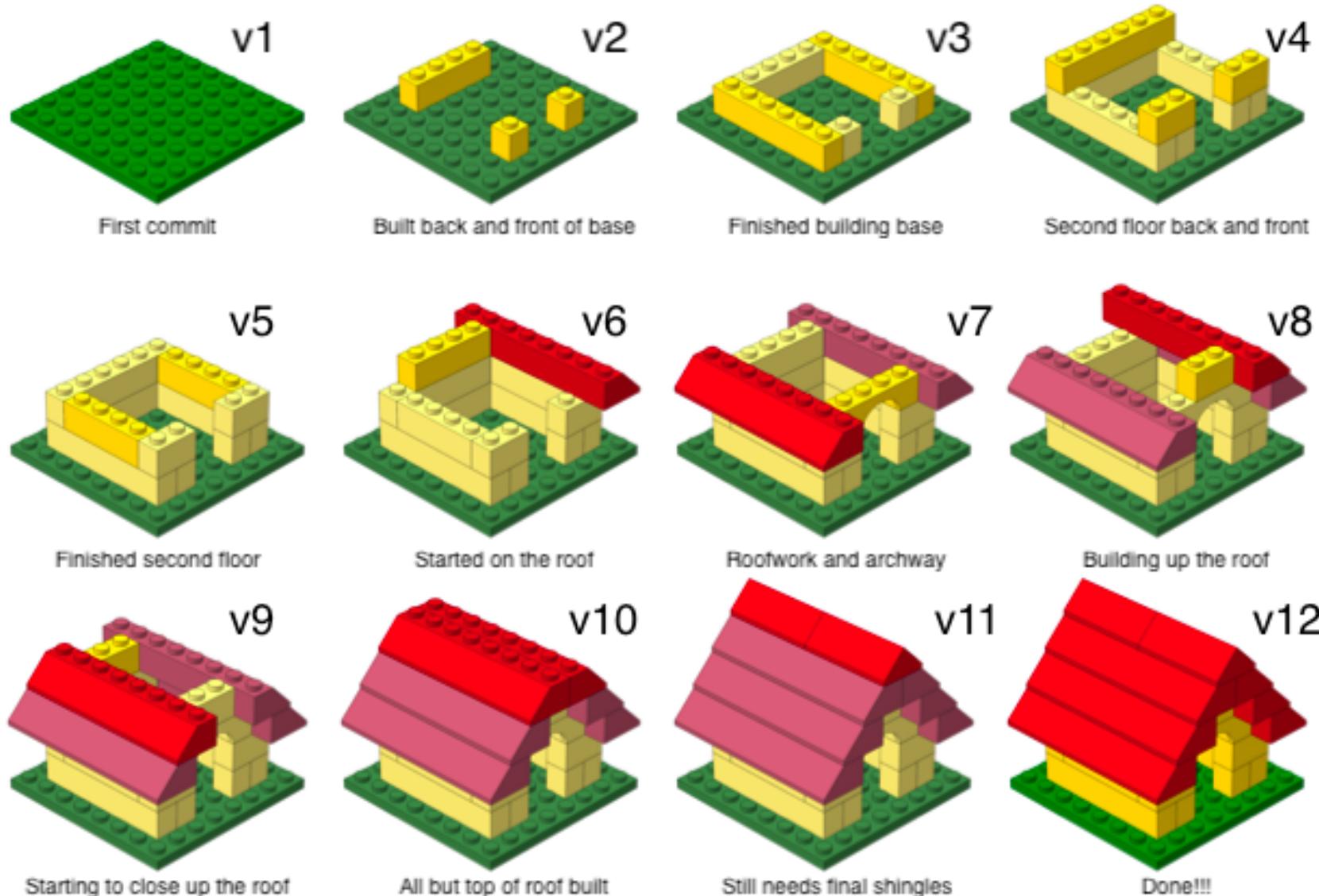


v11

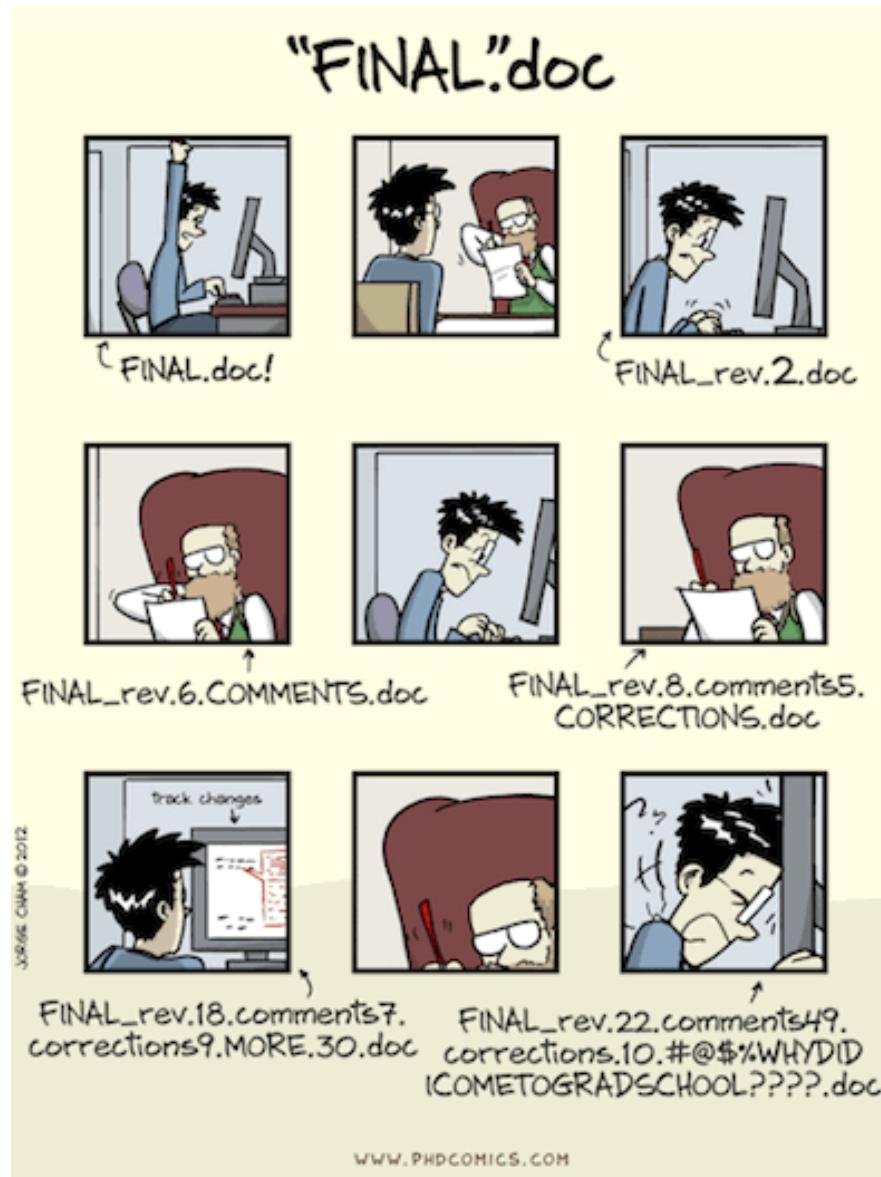


v12

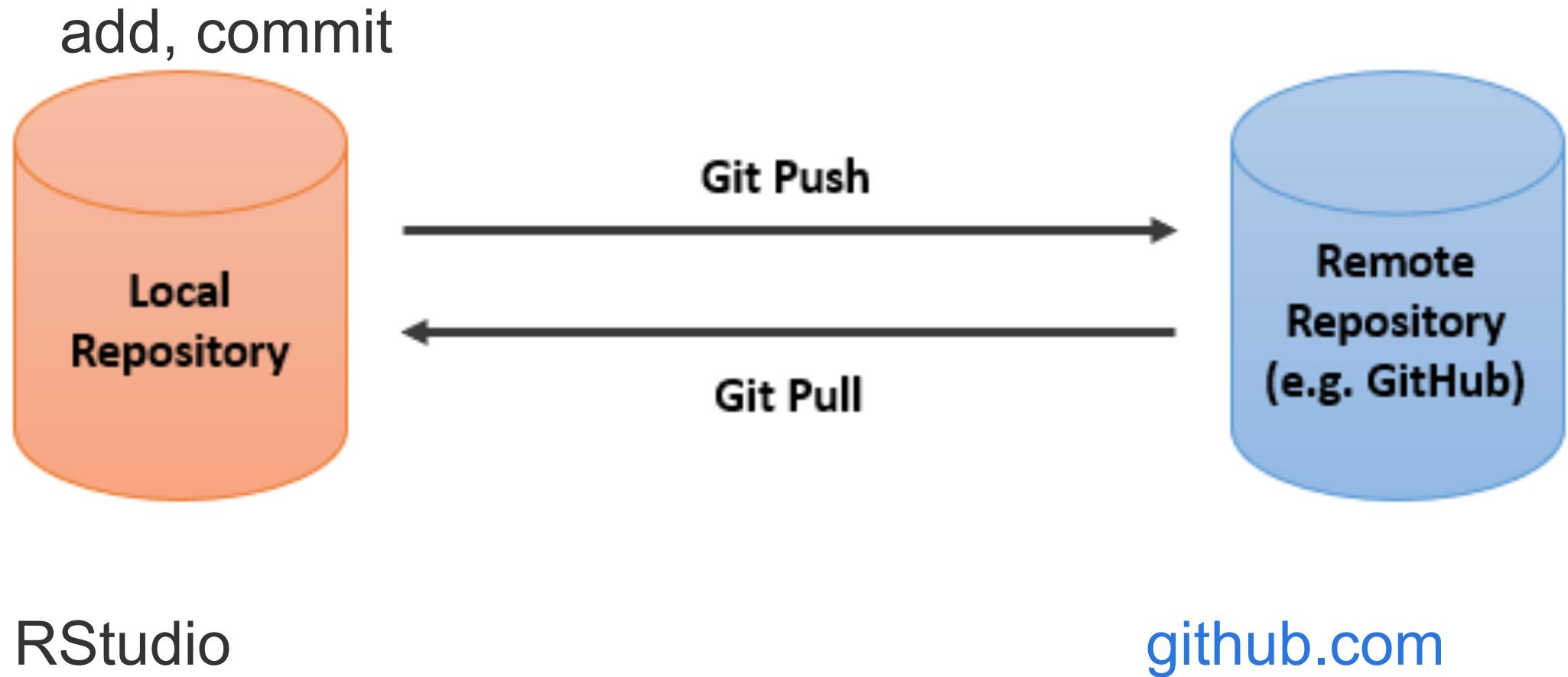
Versionierung



Warum Versionierung



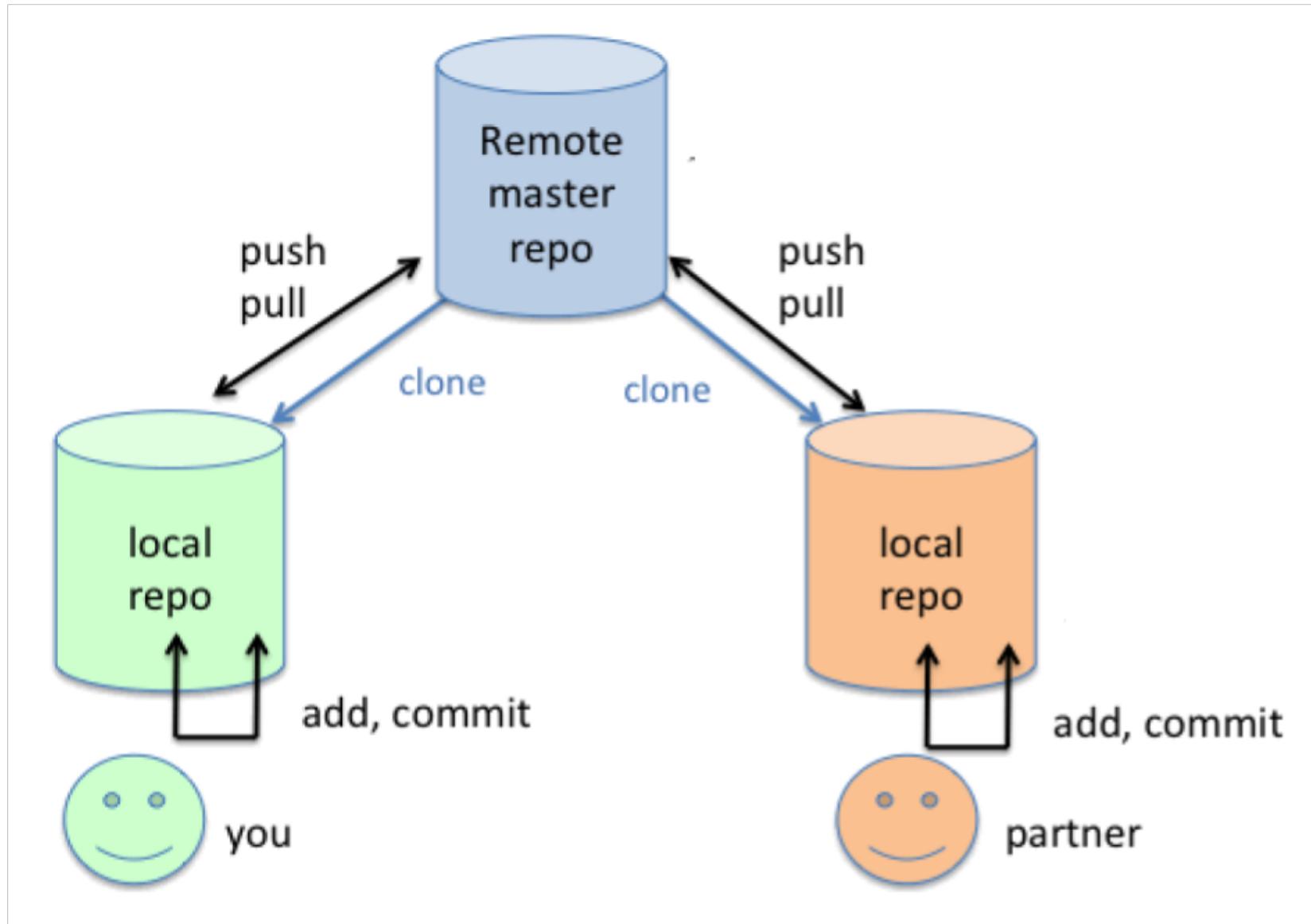
Kollaboration mit sich Selbst



RStudio

github.com

Kollaboration mit anderen



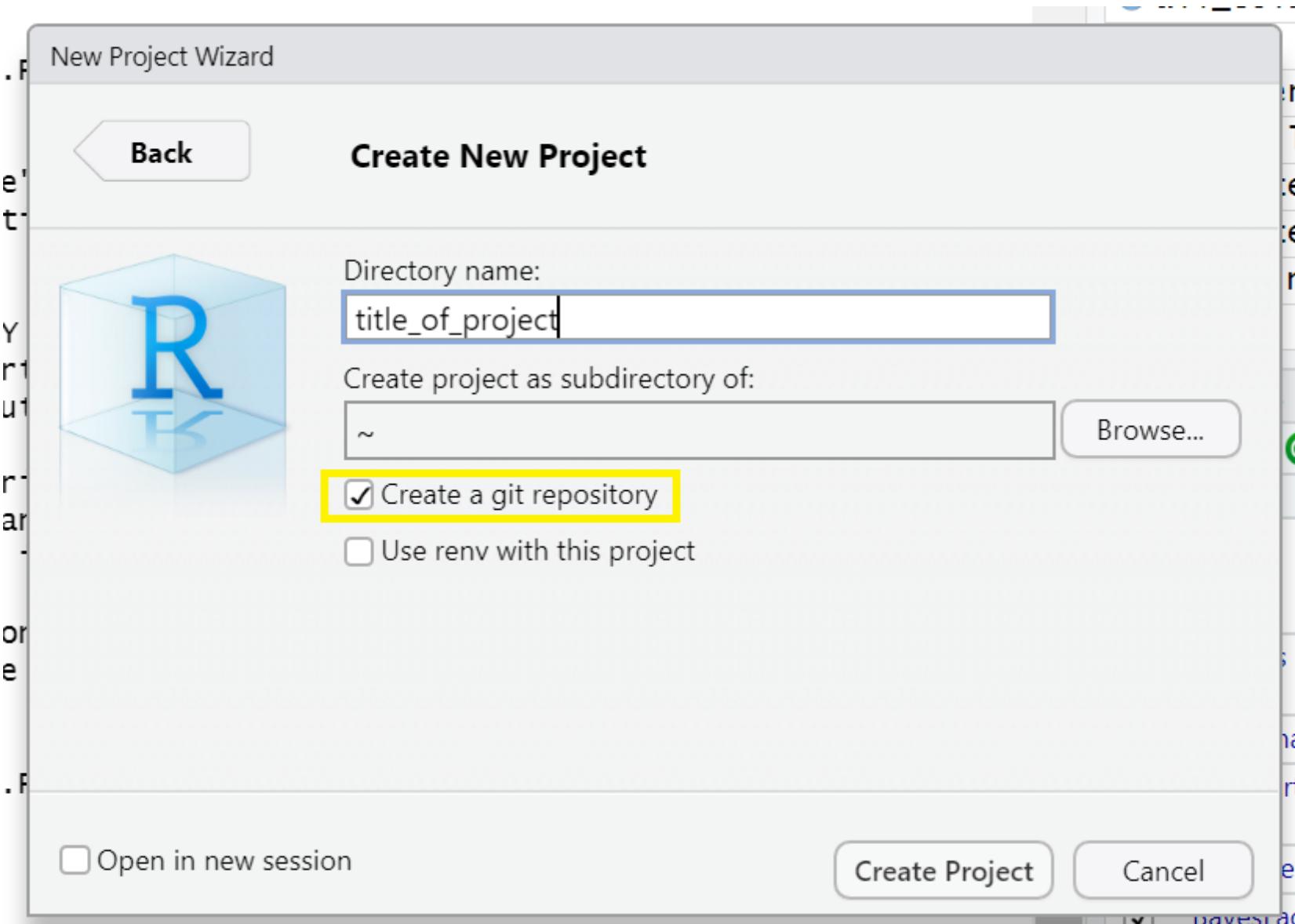
Warum Git und GitLab?



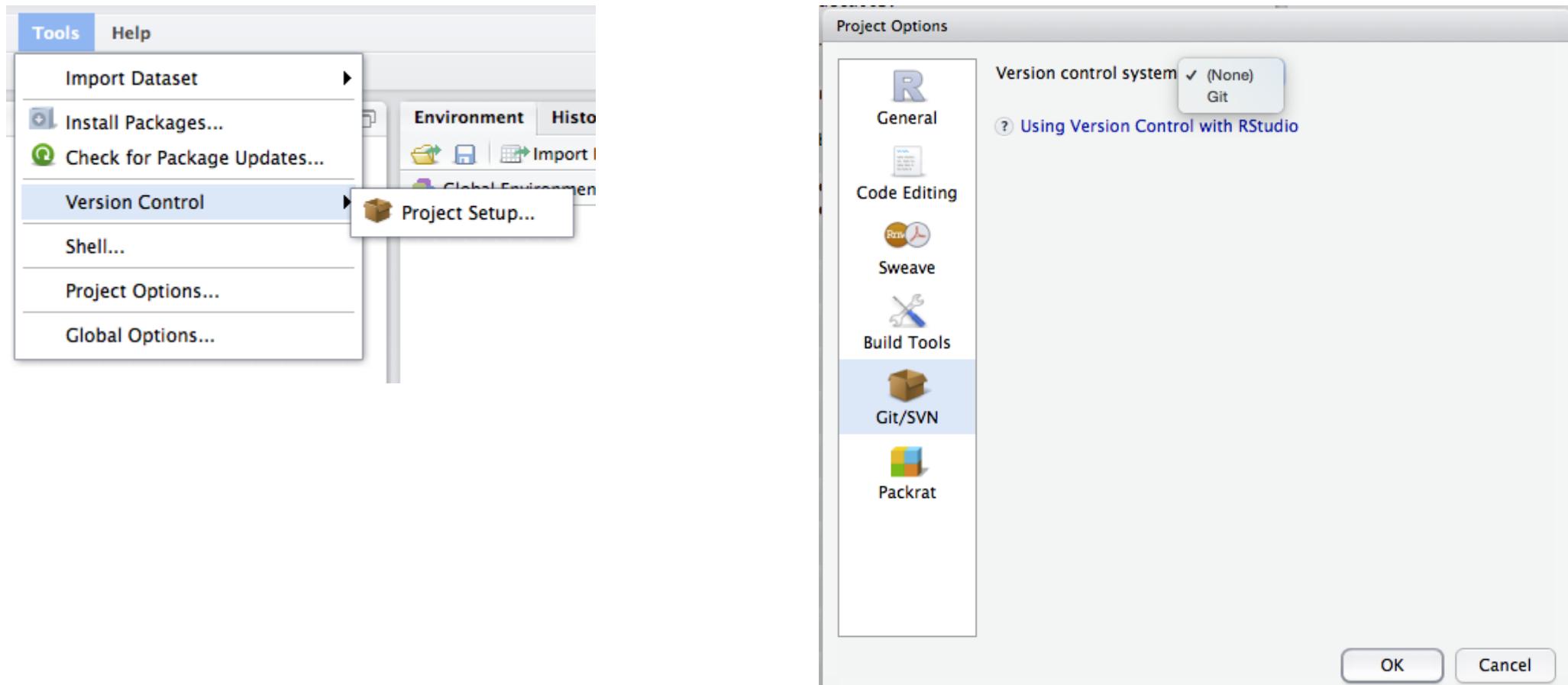
Git und GitLab

- (Nicht ganz) Millionen von Befehlen
- 99% der Zeit: `add`, `commit`, `push`, `pull`
- In RStudio oder *command line (terminal)*
- Happy Git and GitHub for the useR

Neues Projekt

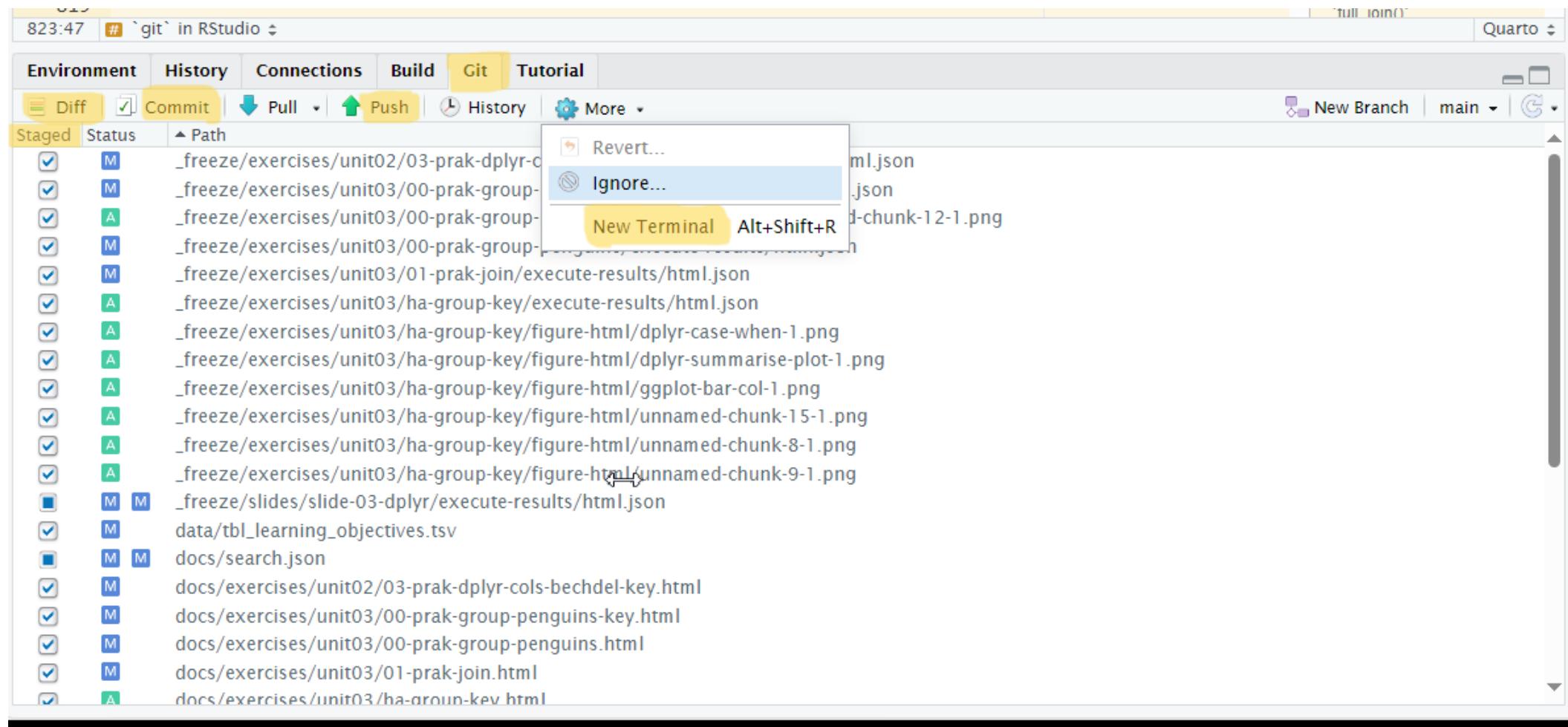


Existierendes Projekt



git in RStudio

Step-by-Step 1, step-by-step 2



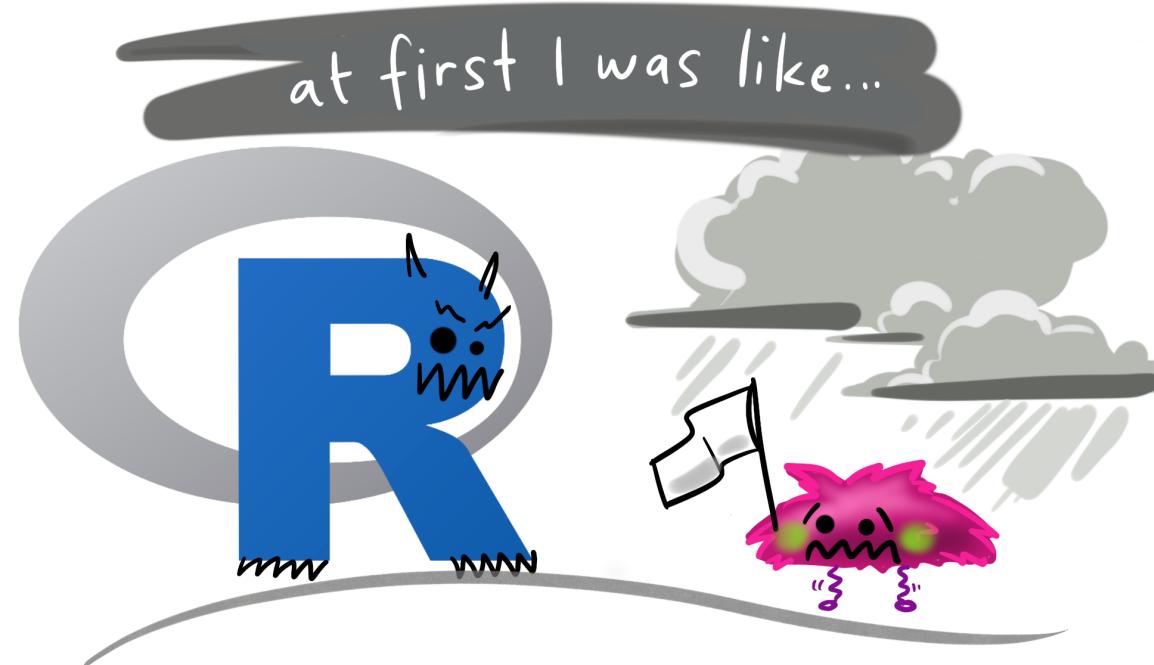


Source: [xkcd](https://xkcd.com/1192/)

Praktikum 06b: git

Quarto-Projekt mit git

20:00
rstatsBL - Data Science mit R



...but now it's like...



rstatsBL - Data Science mit R
Horst '18

Fragen?

- R for Data Science
- DataCamp
- Stackoverflow
- ChatGPT, Perplexity Ai
- Souad Guemghar & Luca Hüslер

Danke! 🍑

Slides created via [revealjs](#) and Quarto.

Access slides as [PDF](#).

All material is licensed under [Creative Commons Attribution Share Alike 4.0 International](#).

Git-Begriffe

- **repository** dein Projektordner
- **commit** ein Schnappschuss deines Repos
- **push** Commits an den *remote* senden
- **pull** Commits vom *remote* erhalten
- **clone** das Repository zum ersten Mal aus dem *remote* abrufen
- **remote** einen Computer, auf dem sich das Repository befindet