

## 02 CSI online aphasia: Typing - Preprocessing

Kirsten Stark

19 August, 2021

@ Marcus: Am besten, du klickst in dem Ordner, den ich dir als .zip-file geschickt habe, die Datei “CSI\_online\_aphasia.Rproj” an. Dann öffnet sich R Studio neu. In dem neuen Fenster unten links unter Files -> Scripts findest du dann dieses Skript hier mit dem Namen “02\_CSI\_online\_aphasia\_typing\_preprocessing.Rmd”). Doppelklick darauf sollte das Skript öffnen. Jetzt kannst du bei jedem grau hinterlegten Feld nacheinander auf das grüne Dreieck drücken und die einzelnen Abschnitte laufen lassen. Dabei am besten immer auf die Kommentare achten, ob du noch etwas anpassen musst =)

Der Code mit # ist jeweils auskommentiert, alles in grauen Boxen oder # davor wird läuft bei Drücken auf Play durch.

### Load packages

```
rm(list = ls())

# install.packages("remotes") # uncomment if installation is needed (only once)
# remotes::install_github("rstudio/renv") # uncomment if installation is needed (only once)
# if file is accessed through the R package, all packages should be installed if
# renv::restore()
# is applied. Otherwise use:
# install.packages("tidyr") # uncomment if installation is needed (only once)
# install.packages("dplyr") # uncomment if installation is needed (only once)
# install.packages("here") # uncomment if installation is needed (only once)
# install.packages("knitr")

library(tidyr)
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union
```

```
options( "encoding" = "UTF-8" )
```

```
# set working directory -- @ Marcus: Am besten, du klickst in dem Ordner, den ich dir als .zip-file ges
# setwd("/Users/kirstenstark/Documents/Research/CSI_online_aphasia")
```

## Load and preprocess data

This input file needs to be entered by hand:

```
# input output main data
type <- c("pretest_patient")
input <- c("data_dissmtukl_2021-08-19_09-30.csv")
#input <- c("data_dissmtukl_2021-08-12_22-27.csv") # can be a vector of several files
# possibly, the CSV file needs to be opened once and saved as CSV UTF-8 (durch Trennzeichen getrennte D
# then all relevant rows of column A need to be selected (e.g. A1:3). Then click Daten --> Text in Spal
# save again as CSV UTF-8
```

```
options( "encoding" = "UTF-8" )
```

```
# output
output <- c("pretest_typing_patient_long.csv", # data file for typewritten data
            "pretest_spoken_patient_long.csv") # data file for spoken data
```

```
# arrays
arrays <- "arrays_umlaut.csv"
```

```
# load data
datafiles <- list()
for(i in 1:length(input)) {
  datafiles[[i]] <- read.csv(here::here("data", "raw", input[i]), sep = ";",
                             na = "")
}
```

```
## Warning in read.table(file = file, header = header, sep = sep, quote
## = quote, : unvollständige letzte Zeile von readTableHeader in '/
## Users/kirstenstark/Documents/Research/CSI_online_aphasia/data/raw/
## data_dissmtukl_2021-08-19_09-30.csv' gefunden
```

```
# perform some transformations on each dataframe
for(i in 1:length(input)) {
  # add type column
  datafiles[[i]]$type <- type[i]
  # delete description column and experimenter's tryout data
  # datafiles[[i]] <- datafiles[[i]][-c(1:2),]
  # add subject id
  datafiles[[i]] <- datafiles[[i]] %>% dplyr::mutate(subject = row_number())
}
```

## Convert to long format, prepare wide dataframe, and bind long and wide dataframes together

First convert all variables with values for each trial, then bind them together. In a next step bind them to the variables that only have one value per participant.

```
for(i in 1:length(input)){  
  
  #-----  
  # Prepare long data frame  
  ## if the file also contains typing data  
  if("online_CSI_typing" %in% datafiles[[i]]$QUESTNNR){  
    ### MAIN TASK - LETTERS  
    # letters of first 1-80 trials of the main experiment  
    df1 <- datafiles[[i]] %>% select('subject', starts_with(c("XT"))) %>%  
      mutate(across(starts_with("XT"), as.character)) %>%  
        pivot_longer(  
          cols = -subject,  
          names_to = c("trial", ".value"),  
          names_pattern = "([~_]+)_(.*)",  
          values_to = "timing") %>%  
        group_by(subject) %>%  
        dplyr::mutate(trial = seq(1,80, by = 1)) %>%  
        setNames(paste0('letters.', names(.)))  
    # letters of last 81-160 trials of the main experiment  
    df2 <- datafiles[[i]] %>% select('subject', starts_with(c("XU"))) %>%  
      mutate(across(starts_with("XU"), as.character)) %>%  
        pivot_longer(  
          cols = -subject,  
          names_to = c("trial", ".value"),  
          names_pattern = "([~_]+)_(.*)",  
          values_to = "timing") %>%  
        group_by(subject) %>%  
        dplyr::mutate(trial = seq(81,160, by = 1)) %>%  
        setNames(paste0('letters.', names(.)))  
    # bind first 80 and last 80 trials together  
    df_letters <- bind_rows(df1, df2) %>%  
      dplyr::rename(subject = letters.subject,  
                    trial = letters.trial) %>%  
      arrange(subject, trial)  
    # delete letter columns from wide data frame  
    datafiles[[i]] <- datafiles[[i]] %>% select(!starts_with(c("XT", "XU")))  
  
    #### MAIN TASK - LETTER TIMING  
    # letter timing of first 1-80 trials of the main experiment  
    df1 <- datafiles[[i]] %>%  
      select('subject', starts_with("TI") &  
            !starts_with("TIME")) %>%  
        pivot_longer(  
          cols = -subject,  
          names_to = c("trial", ".value"),  
          names_pattern = "([~_]+)_(.*)",  
          values_to = "timing") %>%  
  }
```

```

      group_by(subject) %>%
      dplyr::mutate(trial = seq(1,80, by = 1)) %>%
      setNames(paste0('timing.', names(.)))
# letter timing of last 81-160 trials of the main experiment
df2 <- datafiles[[i]] %>% select('subject', starts_with(c("TJ"))) %>%
  pivot_longer(
    cols = -subject,
    names_to = c("trial", ".value"),
    names_pattern = "([~_]+)_(.*)",
    values_to = "timing") %>%
    group_by(subject) %>%
    dplyr::mutate(trial = seq(81,160, by = 1)) %>%
    setNames(paste0('timing.', names(.)))
# bind first 80 and last 80 trials together
df_timing <- bind_rows(df1, df2) %>%
  dplyr::rename(subject = timing.subject,
    trial = timing.trial) %>%
  arrange(subject, trial)
# delete timing columns from wide data frame
datafiles[[i]] <- datafiles[[i]] %>%
  select(!(starts_with(c("TI", "TJ")) & !starts_with("TIME")))

### MAIN TASK - ENTIRE WORDS
# typed words of first 1-80 trials of the main experiment
df1 <- datafiles[[i]] %>%
  select('subject', starts_with("AW")) %>%
  mutate(across(starts_with("AW"), as.character)) %>%
  pivot_longer(
    cols = -subject,
    names_to = c("trial"),
    values_to = "word") %>%
    group_by(subject) %>%
    dplyr::mutate(trial = seq(1,80, by = 1))
# typed words of last 1-80 trials of the main experiment
df2 <- datafiles[[i]] %>%
  select('subject', starts_with("AV")) %>%
  mutate(across(starts_with("AV"), as.character)) %>%
  pivot_longer(
    cols = -subject,
    names_to = c("trial"),
    values_to = "word") %>%
    group_by(subject) %>%
    dplyr::mutate(trial = seq(81,160, by = 1))
# bind first 80 and last 80 trials together
df_words <- bind_rows(df1, df2) %>%
  arrange(subject, trial)
# delete word columns from wide data frame
datafiles[[i]] <- datafiles[[i]] %>%
  select(!starts_with(c("AW", "AV")))

### FAMILIARIZATION WORDS
### (only to control that ppt payed attention to the task)

```

```

# typed words of first 1-80 trials of the familiarization
df1 <- datafiles[[i]] %>%
  select('subject', starts_with("FA03")) %>%
  pivot_longer(
    cols = -subject,
    names_to = c("trial"),
    values_to = "fam_typed") %>%
  group_by(subject) %>%
  dplyr::mutate(trial = seq(1,80, by = 1))
# typed words of last 1-80 trials of the main experiment
df2 <- datafiles[[i]] %>%
  select('subject', starts_with("FA04")) %>%
  pivot_longer(
    cols = -subject,
    names_to = c("trial"),
    values_to = "fam_typed") %>%
  group_by(subject) %>%
  dplyr::mutate(trial = seq(81,160, by = 1))
# bind first 80 and last 80 trials together
df_fam <- bind_rows(df1, df2) %>%
  arrange(subject, trial)

# delete word columns from wide data frame
datafiles[[i]] <- datafiles[[i]] %>%
  select(!starts_with(c("FA03", "FA04")))
}

## if the file also contains spoken data
if("online_CSI_spoken" %in% datafiles[[i]]$QUESTNNR){
  ### AUDIO FILES
  # Audio files of first 1-80 trials
  df1 <- datafiles[[i]] %>%
    select('subject', starts_with("AR")&contains("x")) %>%
    pivot_longer(
      cols = -subject,
      names_to = c("trial"),
      values_to = "audio") %>%
      group_by(subject) %>%
      dplyr::mutate(trial = seq(1,80, by = 1))
  # Audio files of last 1-80 trials
  df2 <- datafiles[[i]] %>%
    select('subject', starts_with("AU")&contains("x")) %>%
    pivot_longer(
      cols = -subject,
      names_to = c("trial"),
      values_to = "audio") %>%
      group_by(subject) %>%
      dplyr::mutate(trial = seq(81,160, by = 1))
  # bind first 80 and last 80 trials together
  df_audio <- bind_rows(df1, df2) %>%
    arrange(subject, trial)
  # delete audiofile columns from wide data frame

```

```

    datafiles[[i]] <- datafiles[[i]] %>%
      select(!starts_with(c("AR", "AU")))
  }

### Bind the four/five trial-wise dataframes together
if("online_CSI_spoken" %in% datafiles[[i]]$QUESTNNR &
    "online_CSI_typing" %in% datafiles[[i]]$QUESTNNR){
  df_main <- merge(df_fam, df_words, by = c("subject", "trial"))
  df_main <- merge(df_main, df_letters, by= c("subject", "trial"))
  df_main <- merge(df_main, df_timing, by= c("subject", "trial"))
  df_main <- merge(df_main, df_audio, by= c("subject", "trial")) %>%
    arrange(subject, trial)
} else if ("online_CSI_typing" %in% datafiles[[i]]$QUESTNNR){
  df_main <- merge(df_fam, df_words, by = c("subject", "trial"))
  df_main <- merge(df_main, df_letters, by= c("subject", "trial"))
  df_main <- merge(df_main, df_timing, by= c("subject", "trial")) %>%
    arrange(subject, trial)
} else {
  df_main <- df_audio %>% arrange(subject, trial)
}

#-----
# Adapt wide data frame with info that is assessed only once

# for control reasons: calculate time sum by hand:
# sum dwell times for each page
datafiles[[i]] <- datafiles[[i]] %>%
  mutate_at(vars(contains("TIME0")), as.numeric)
datafiles[[i]] <- datafiles[[i]] %>% rowwise() %>%
  dplyr::mutate(timetotal = rowSums(across(starts_with("TIME0")))/60)

# delete columns with info we don't need
datafiles[[i]] <- datafiles[[i]] %>%
  select(-c(CASE, MODE, STARTED, SD22_PRV,
            SD22_BID, SD22_BVS,
            SD28_01, SD29_01,
            FA02_01, PT01, PT02, PT03, PT04,
            LASTDATA, SD19, SD19_01, SD19_02, SD19_03,
            MISSING, MISSREL
            ))

# delete columns that contain only NAs
datafiles[[i]]<- datafiles[[i]] %>% select_if(~sum(!is.na(.)) > 0)

# delete practice audio files
if("online_CSI_spoken" %in% datafiles[[i]]$QUESTNNR){
  datafiles[[i]]<- datafiles[[i]] %>% select(!starts_with("PA"))
}

# give columns more recognizable names
if (type[i] == "pretest_patient") {
  datafiles[[i]] <- datafiles[[i]] %>%
    dplyr::rename(array_no = AY01_01, gender = SD01, age = SD02_01,

```

```

        language.test = SD20, language = SD21,
        os_system = SD22_OS, browser_automatic = SD22_BNM,
        system_format = SD22_FmF,
        fingers_l = SD25, fingers_r = SD26,
        handedness = SD27, operator_system = SD30,
        system = SD32, browser = SD31,
        typing_experience = SD33,
        keyboard_type = KB01,
        time_wo_outlier = TIME_SUM,
        screen_width = SD22_ScW, screen_height = SD22_ScH,
        questionnaire_width = SD22_QnW, array=OR01_01)
} else if (type[i] == "main") {
  datafiles[[i]] <- datafiles[[i]] %>%
  dplyr::rename(array_no = AY01_01, gender = SD01, age = SD02_01,
    language.test = SD20, language = SD21,
    os_system = SD22_OS, browser_automatic = SD22_BNM,
    system_format = SD22_FmF, prolificid = SD24_01,
    fingers_l = SD25, fingers_r = SD26,
    handedness = SD27, operator_system = SD30,
    system = SD32, browser = SD31,
    browser_other = SD31_07,
    keyboard_type = KB01,
    comments = IM01_01, time_wo_outlier = TIME_SUM,
    screen_width = SD22_ScW, screen_height = SD22_ScH,
    questionnaire_width = SD22_QnW)
} else if (type[i] == "replacement") {
  datafiles[[i]] <- datafiles[[i]] %>%
  dplyr::rename(array_no = AY01_01, gender = SD01, age = SD02_01,
    language.test = SD20, language = SD21,
    os_system = SD22_OS, browser_automatic = SD22_BNM,
    system_format = SD22_FmF, prolificid = SD24_01,
    fingers_l = SD25, fingers_r = SD26,
    handedness = SD27, operator_system = SD30,
    system = SD32, browser = SD31,
    operator_system_other = SD30_07,
    keyboard_type = KB01,
    comments = IM01_01, time_wo_outlier = TIME_SUM,
    screen_width = SD22_ScW, screen_height = SD22_ScH,
    questionnaire_width = SD22_QnW,
    array=OR01_01)
}

#-----
# Bind long and wide data frame together
# Repeat each subjects' rows 160 times (no of trials)
datafiles[[i]] <- datafiles[[i]] %>% slice(rep(seq_len(n()), 160))

# Add trial number to wide data frame
datafiles[[i]]$trial <- rep(1:160, times = max(datafiles[[i]]$subject))

# bind wide and long info together
datafiles[[i]] <- datafiles[[i]] %>%
  left_join(df_main, by = c("subject", "trial")) %>%

```

```

relocate(subject, trial)

#-----
# Convert numeric variables from string to integer:
#str(df)
if(type[i] == "pretest"){
  datafiles[[i]] <- datafiles[[i]] %>%
  mutate_at(vars(!c("prolificid",
                    "word", "type",
                    "array_no", contains("KB0"),
                    contains("TT0"), contains("fam_typed"),
                    contains("letters"))), as.numeric)
}else if(type[i] == "main") {
  datafiles[[i]] <- datafiles[[i]] %>%
  mutate_at(vars(!c("prolificid", "browser_other",
                    "word", "comments", "type",
                    "array_no", "TIME_RSI", contains("KB0"),
                    contains("TT0"), contains("fam_typed"),
                    contains("letters"))), as.numeric)
} else if(type[i] == "replacement") {
  datafiles[[i]] <- datafiles[[i]] %>%
  mutate_at(vars(!c("prolificid", "operator_system_other",
                    "word", "comments", "type",
                    "array_no", contains("KB0"),
                    contains("TT0"), contains("fam_typed"),
                    contains("letters"))), as.numeric)
}
}

```

## Roughly check participants' adherence to the experiment

Did all participants finish the experiment?

```

# did all participants finish the experiment?
for(i in 1:length(input)) {
  print(paste(type[i], ":"))
  print("Experiment completed?")
  print(table(datafiles[[i]]$FINISHED)/160)
  print("What was the last experimental page reached?")
  print(table(datafiles[[i]]$LASTPAGE)/160)
}

```

```

## [1] "pretest_patient :"
## [1] "Experiment completed?"
##
## 1
## 2
## [1] "What was the last experimental page reached?"
##
## 25 33
## 1 1

```



Exclude participants who did not finish the experiment for approval:

```
# for (i in 1:length(input)) {
#   exclude <- datafiles[[i]]$prolificid[datafiles[[i]]$FINISHED == "0" &
#                                     datafiles[[i]]$LASTPAGE != "30"]
#   exclude_id <- datafiles[[i]]$subject[datafiles[[i]]$FINISHED == "0" &
#                                     datafiles[[i]]$LASTPAGE != "30"]
#
#   # Exclude from data frame
#   datafiles[[i]] <- datafiles[[i]][!datafiles[[i]]$subject %in% exclude_id,]
#
#   # update subject IDs
#   datafiles[[i]]$subject <-
#     rep(1:nlevels(as.factor(datafiles[[i]]$subject)), each = 160)
#
#   # All the other data look fine, so we can include all others:
#   include <- as.data.frame(
#     datafiles[[i]]$prolificid[datafiles[[i]]$FINISHED == "1" &
#                               datafiles[[i]]$LASTPAGE == "33"])
#
#   # delete prolific ids to anonymize data frame
#   datafiles[[i]] <- datafiles[[i]] %>%
#     dplyr::select(-prolificid)
# }
#
```

## Add array (actual stimuli) for each participant

```
# load arrays
arrays <- read.csv(here::here("data", "supplementary_info", arrays),
  sep = ";", na = "NA")

for (i in 1:length(input)) {
  # convert array column in datafiles
  datafiles[[i]]$array_no <- as.numeric(stringr::str_remove(
    datafiles[[i]]$array_no, "Array"))
  # bind arrays to data frame
  datafiles[[i]]$item <- NA
  for(j in 1:nlevels(as.factor(datafiles[[i]]$subject))){
    array <- arrays[, unique(datafiles[[i]]$array[
      datafiles[[i]]$subject == j])]
    datafiles[[i]]$item[datafiles[[i]]$subject == j] <- array
  }

  # add category columns
  datafiles[[i]] <- datafiles[[i]] %>% group_by(subject, item) %>%
    dplyr::mutate(category = arrays$categorie[arrays$item == item]) %>%
    dplyr::mutate(supercategory =
      arrays$supercategorie[arrays$item == item] )

  # add position number
}
```

```

datafiles[[i]] <- datafiles[[i]] %>% group_by(subject, category) %>%
  add_count() %>%
  dplyr::mutate(PosOr = seq(1:n)) %>% select(-n)
# add ordinal position for fillers
table(datafiles[[i]]$PosOr)
count <- 1
for (j in 1:nrow(datafiles[[i]])) {
  if(datafiles[[i]]$category [j] == "Filler") {
    datafiles[[i]]$PosOr[j] <- count
    count <- count+1
  }
  if(count == 6) { count <- 1}
}
table(datafiles[[i]]$PosOr)
}

```

## Bind the dataframes together

```
df <- bind_rows(datafiles)
```

## Split dataframes into spoken und typed data

```

df_typing <- df %>% filter(QUESTNNR=="online_CSI_typing") %>%
  select_if(~sum(!is.na(.)) > 0)
df_spoken <- df %>% filter(QUESTNNR=="online_CSI_spoken") %>%
  select_if(~sum(!is.na(.)) > 0)

```

## Fix single participants

In the typing experiment, the typed words of participant 1 were not saved (reason unknown). We fix this by putting all the single letter columns together. For some reason, space or enter were also saved as NANA. We will fix that later

```
df_typing <- df_typing %>% unite(word,starts_with("letter"),sep="", remove=FALSE, na.rm=TRUE)
```

## Export prepared data frame

```

if("online_CSI_typing" %in% unique(df$QUESTNNR)) {
  write.csv(df_typing, here::here("data", "transient_data_files", output[1]), row.names = FALSE)
}
if("online_CSI_spoken" %in% unique(df$QUESTNNR)) {
  write.csv(df_spoken, here::here("data", "transient_data_files", output[2]), row.names = FALSE)
}

```