

# Practical work: Editing raw data from sniffer



Oscar González-Recio and Coralia Manzanilla-Pech





# Overview: First part-Connecting

- 1. Starting with Jupyter Notebook in Annuna (WUR)
- 2. Opening Github to see the files
- 3. Cloning files from Github to Jupyter Notebook
- 4. Ready for practicals



#### Preliminars: How to connect to the WUR server

https://notebook.anunna.wur.nl





#### Github

Oscar's: <a href="https://ogrecio.github.io/RelivestockMethaneCourse/">https://ogrecio.github.io/RelivestockMethaneCourse/</a>

Coralia's: <a href="https://github.com/cmanzanillap">https://github.com/cmanzanillap</a>

Ester's: <a href="https://github.com/estermt/sniffer Data management">https://github.com/estermt/sniffer Data management</a>



#### Environments







#### Practical environment



https://notebook.anunna.wur.nl

test.ipynb

# Server Options Jupyterhub on Anunna

Select a location for your serve	:	
Cluster - Re-livestock		
	Start	



# How it looks





# How to clone files/folders from Github to Jupyter

- 1. Open Jupyter Notebook online in your web browser.
- 2. Click on the "New" button in the top right corner of the screen.
- $3.\;\;$  Select "Terminal" from the dropdown menu.
- 4. In the terminal window, type the following command:

git clone https://github.com/ogrecio/RelivestockMethaneCourse



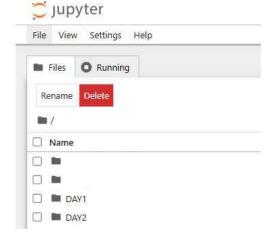
#### How it looks

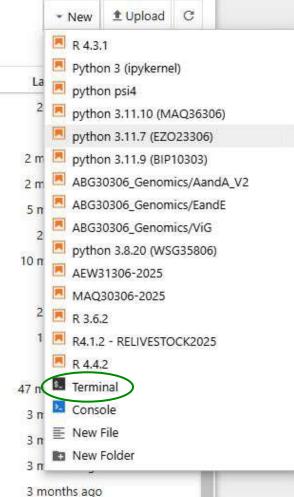
#### RelivestockMetaneCourse/Day1

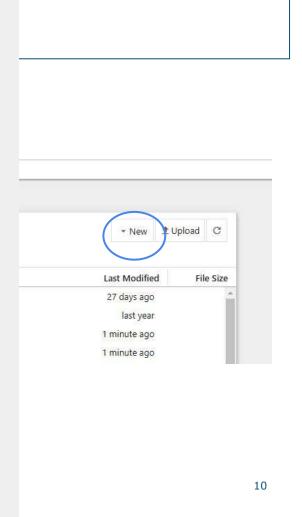




#### How it looks

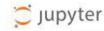








# How it looks Terminal



```
File View Settings Help
manza003@node257:~$ git clone https://github.com/cmanzanillap/DAY1
Cloning into 'DAY1' ...
remote: Enumerating objects: 34, done.
remote: Counting objects: 100% (34/34), done.
remote: Compressing objects: 100% (34/34), done.
remote: Total 34 (delta 14), reused 0 (delta 0), pack-reused 0 (from 0)
                                                                                       type cd Day1
Unpacking objects: 100% (34/34), 5.49 MiB | 4.59 MiB/s, done.
manza003@node257:~$ git clone https://github.com/cmanzanillap/DAY2
Cloning into 'DAY2' ...
remote: Enumerating objects: 35, done.
remote: Counting objects: 100% (35/35), done.
remote: Compressing objects: 100% (35/35), done.
remote: Total 35 (delta 18), reused 0 (delta 0), pack-reused 0 (from 0)
Unpacking objects: 100% (35/35), 376.44 KiB | 1.76 MiB/s, done.
manza003@node257:~$ cd DAY1
manza003@node257:~/DAY1$ ls
                                                 herd sniffer101 Loggy 42s FD1.txt.gz output.txt
Rscript.ipynb herd robot101.csv
                                                                                                        test.ipynb
              herd sniffer101 Loggy 42s FD1.txt input raw.txt
                                                                                      script slurm2.sh
bash, sh
manza003@node257:~/DAY1$
```



#### useful commands Linux

```
Is (to see files) Is -d */ (to see directories)
less (to see the files inside) less -S (useful for genotypes) a
(to quit from less)
head (head of the file)
wc -l file (number of rows)
awk '{print NF; exit}' file (number of columns)
cd folder (change directory) cd .. (comeback to the previous level)
rm -rf directory or rm file
```



# Files in RelivestockMethaneCourse/Day1 folder

- 1. Database to read in R: input\_raw.txt
- 2. Jupyter R Notebook: Rscript.ipynb, merge\_sniffer\_testday.ipynb
- Databases to run Java program: herd\_robot101.csv herd\_sniffer101\_Loggy\_42s\_FD1.txt
- 4. Script to run in a bash file: bash.sh



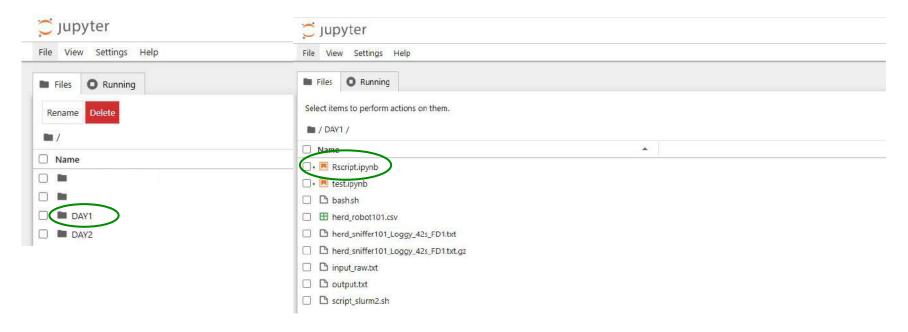
# Overview: Second part- Practical in R

- 1. Installing packages
- 2. Read data input\_raw.txt
- 3. Get familiar with the data (check number records, number cows, duration, min max etc)
- 4. Create event (visit), plotting events to see how data looks
- 5. Background approaches
- 6. Calculate different phenotypes
- 7. Calculate correlations between them



#### How it looks

#### RelivestockMetaneCourse/Day1





# Raw CH<sub>4</sub> data from sniffer

```
Jupyter output.txt Last Checkpoint: 2 minutes ago
File Edit View Settings Help
    1 cow, date, ch4, co2
    2 59,2022/18/02 00:00:01,0.042,0.512
    3 59,2022/18/02 00:00:02,0.042,0.52
    4 59,2022/18/02 00:00:03,0.041,0.521
    5 59,2022/18/02 00:00:04,0.041,0.526
    6 59,2022/18/02 00:00:05,0.041,0.527
    7 59,2022/18/02 00:00:06,0.041,0.529
    8 59,2022/18/02 00:00:07,0.041,0.533
    9 59,2022/18/02 00:00:08,0.041,0.534
   10 59,2022/18/02 00:00:09,0.04,0.532
   11 59,2022/18/02 00:00:10,0.04,0.534
   12 59,2022/18/02 00:00:11,0.04,0.538
   13 59,2022/18/02 00:00:12,0.04,0.542
   14 59,2022/18/02 00:00:13,0.04,0.542
   15 59,2022/18/02 00:00:14,0.04,0.544
   16 59,2022/18/02 00:00:15,0.04,0.545
   17 59,2022/18/02 00:00:16,0.04,0.549
   18 59,2022/18/02 00:00:17,0.04,0.551
   19 59,2022/18/02 00:00:18,0.04,0.549
   20 59,2022/18/02 00:00:19,0.04,0.549
   21 59,2022/18/02 00:00:20,0.04,0.547
   22 59,2022/18/02 00:00:21,0.041,0.549
   23 59,2022/18/02 00:00:22,0.042,0.549
```

# columns:

date ch4 co2

# Call the following packages





# Read the output.txt file

```
rawCH4 <- read.table("input_raw.txt", sep = ",", stringsAsFactors = F, header = T)
colnames(rawCH4)
head(rawCH4)
dim(rawCH4)
summary(rawCH4)</pre>
```



# Multiplicate to get ppm & format date and time

```
# Multiply to get ppm
rawCH4ch4ppm < -rawCH4ch4 * 10000
rawCH4co2ppm < -rawCH4co2 * 10000

# Parse datetime
rawCH4$datetime <- as.POSIXct(rawCH4$date, format = "%Y/%d/%m %H:%M:%S")
rawCH4$date <- as.Date.character(rawCH4$date, "%Y/%d/%m")
head(rawCH4)</pre>
```



# Create an event per animal

```
#create an event per entrance to the AMS per day per animal
event_ch4 <- rawCH4 %>%
  select(ch4ppm, co2ppm, cow, date, datetime) %>%
  mutate(eventID = 1 + cumsum(cow != lag(cow, default = first(cow))))
hist(event_ch4$eventID)
head(event_ch4)
```



#### Get familiar with the data

```
num_cows <- event_ch4 %>%
  summarise(unique_cows = n_distinct(cow))
print(num_cows)

cow_counts <- rawCH4 %>%
  group_by(cow) %>%
  summarise(count = n())
head(cow_counts)
```



# Calculating duration of the visit

```
######### calculating duration of the visit #############
event_duration <- event_ch4 %>%
 group by(eventID) %>%
  summarize(
   start time = min(datetime),
   end time = max(datetime),
   lenght = as.numeric(difftime(max(datetime), min(datetime), units = "secs"))
head(event duration)
summary(event duration)
hist(event duration$lenght)
```



# Calculating number of visits per cow

```
#Calculating how many events (visit) per cow

event_counts <- event_ch4 %>%
   group_by(cow) %>%
   summarise(unique_events = n_distinct(eventID))

# View the result
print(event_counts)
```



# Calculating background Approach 1

```
event_3low <- event_ch4 %>%
  group_by(eventID) %>%
  arrange(ch4) %>%
  slice_head(n = 3) %>%
  summarise(avg_ch4 = mean(ch4, na.rm = TRUE))
head(event_3low)
summary(event_3low)
```



# Calculating background Approach 2

```
########### calculating the 0.001 quantile per visit to use it as background ##########
event_quant <- event_ch4 %>%
  group_by(eventID) %>%
  summarise(ch4_quantile_0_001 = quantile(ch4, probs = 0.001, na.rm = TRUE))
head(event_quant)
summary(event_quant)
```



# Calculating average concentration for CH4 - CO2



# Now for the cut visit length

```
######### calculating the mean for ch4 only for 60-300 sec per visit ####### losing one event
event cut <- event ch4 %>%
 # Trim any spaces from the TimeStamp column
 mutate(datetime = trimws(datetime)) %>%
 # Convert TimeStamp to datetime object
 mutate(datetime = ymd hms(datetime)) %>%
 group by(eventID) %>%
 # Calculate the time difference in seconds from the first TimeStamp per eventID
 mutate(time diff sec = as.numeric(difftime(datetime, min(datetime), units = "secs"))) %>%
 # Filter to keep only rows where the time difference is between 60 and 300 seconds
 filter(time diff sec >= 60 & time_diff_sec <= 300) %>%
 # Optionally, summarize the data (for example, calculating means of ch4 and co2)
 summarise(
   mean ch4 = mean(ch4, na.rm = TRUE),
   mean co2 = mean(co2, na.rm = TRUE)
head(event cut)
summary(event cut)
dim(event cut)
```



# Plotting events to see patterns

```
#Printing cows numbers
cow_ids<- unique(event_ch4$cow)
print(cow_ids)</pre>
```



# Function: Detecting peaks

```
# Function to detect peaks
detectPeaks <- function(ch4, windowSize = 5, threshold = 0.0005) {
 peaksIni <- c()
 peaksFin <- c()
 stage <- FALSE
 peak_decreasing <- FALSE
 for (i in 1:(length(ch4) - windowSize)) {
   tempW <- ch4[i:(i + windowSize - 1)]
   tempX <- 1:windowSize
   mean tempW <- mean(tempW, na.rm = TRUE)
   mean tempX <- mean(tempX, na.rm = TRUE)
   cov_xy <- sum((tempW - mean_tempW) * (tempX - mean_tempX), na.rm = TRUE)
   var x <- sum((tempX - mean_tempX) * (tempX - mean_tempX), na.rm = TRUE)</pre>
   pendiente <- cov xy / var x
   if (pendiente > threshold && !stage) {
     peaksIni <- c(peaksIni, i)
     stage <- TRUE
     peak decreasing <- FALSE
   } else if (stage && pendiente < -threshold) {
     peak decreasing <- TRUE
   } else if (stage && peak decreasing && pendiente > -threshold / 2) {
     peaksFin <- c(peaksFin, i)
     stage <- FALSE
     peak decreasing <- FALSE
   } else if (length(peaksIni) == 0 && pendiente > threshold / 2) {
     # Handle case where no peaks have been detected yet
   } else if (i == (length(ch4) - windowSize) && peak decreasing) {
     peaksFin <- c(peaksFin, i)
 list(peaksIni = peaksIni, peaksFin = peaksFin)
```

# Function: Detecting peaks

```
# Function to detect peaks
detectPeaks <- function(ch4, windowSize = 5, threshold = 0.0005) {
  peaksIni <- c()
 peaksFin <- c()
  stage <- FALSE
  peak_decreasing <- FALSE
  for (i in 1:(length(ch4) - windowSize)) {
    tempW <- ch4[i:(i + windowSize - 1)]
    tempX <- 1:windowSize
    mean tempW <- mean(tempW, na.rm = TRUE)
    mean tempX <- mean(tempX, na.rm = TRUE)
    cov_xy <- sum((tempW - mean_tempW) * (tempX - mean_tempX), na.rm = TRUE)
    var x <- sum((tempX - mean tempX) * (tempX - mean tempX), na.rm = TRUE)</pre>
    pendiente <- cov xy / var x
   if (pendiente > threshold && !stage) {
      peaksIni <- c(peaksIni, i)
      stage <- TRUE
      peak decreasing <- FALSE
    } else if (stage && pendiente < -threshold) {
      peak decreasing <- TRUE
   } else if (stage && peak decreasing && pendiente > -threshold / 2) {
      peaksFin <- c(peaksFin, i)
      stage <- FALSE
      peak decreasing <- FALSE
   } else if (length(peaksIni) == 0 && pendiente > threshold / 2) {
      # Handle case where no peaks have been detected yet
   } else if (i == (length(ch4) - windowSize) && peak_decreasing) {
      peaksFin <- c(peaksFin, i)
 list(peaksIni = peaksIni, peaksFin = peaksFin)
```

```
# Apply the peak detection to the entire database
peak_det <- event_ch4 %>%
   group_by(eventID) %>%
   summarise(num_peaks = length(detectPeaks(ch4)$peaksIni))
# Print the results
head(peak_det)
summary(peak_det)
```

# Exercise: calculate number of peaks per minute



# Function: Detect peaks and calc sum2maxpeaks

```
# Function to detect peaks and calculate sum2maxpeaks
detectPeaksAndSum <- function(ch4, windowSize = 5, threshold = 0.0005) {
  peaksIni <- c()
  peaksFin <- c()
  stage <- FALSE
  peak decreasing <- FALSE
  for (i in 1:(length(ch4) - windowSize)) {
    tempW <- ch4[i:(i + window5ize - 1)]
    tempX <- 1:windowSire
    mean templi <- mean(templi, na.rm = TRUE)
    mean tempX <- mean(tempX, na.rm = TRUE)
    cov xy <- sum[(tempW - mean tempW) * (tempX - mean tempX), na.rm - TRUE]
    var x <- sum((tempX - mean tempX) * (tempX - mean tempX), na.rm = TRUE)
    pendiente <- cov_xy / var_x
    if (pendiente > threshold && !stage) {
     peaksIni <- c(peaksIni, i)
     stage <- TRUE
     peak decreasing <- FALSE
    } else if (stage && pendiente < -threshold) [
     peak_decreasing <- TRUE
    } else if (stage && peak decreasing && pendiente > -threshold / 2) {
     peaksFin <- c(peaksFin, 1)
     stage <- FALSE
     peak decreasing <- FALSE
    } else if (length(peaksIni) == 0 && pendiente > threshold / 2) {
     # Handle case where no peaks have been detected yet
    } else if (i == (length(ch4) - windowSize) && peak decreasing) {
     peaksFin <- c(peaksFin, i)
  sum2maxpeaks <- 8
  if (length(peaksIni) > @ && length(peaksFin) > @) (
   for (j in 1:length(peaksIni)) {
     start <- peaksIni[i]
     end <- ifelse() <= length(peaksFin), peaksFin[j], NA)
     if (!is.na(end) && end > start) {
       peak values <- ch4[start:end]
       max values <- sort(peak values, decreasing = TRUE)[1:2]
       avg max values <- mean(max values)
        sum2maxpeaks <- sum2maxpeaks + avg max values
```

```
# Apply the peak detection to the entire database
peak phen <- event ch4 %>%
 group by(eventID) %>%
 summarise(
   num peaks = {
     result <- detectPeaksAndSum(ch4)
     length(result$peaksIni)
    sum2maxpeaks = {
     result <- detectPeaksAndSum(ch4)
     result$sum2maxpeaks
# Print the results
head(peak phen)
```

# Adding area under the curve

```
sum2maxpeaks <- 0
 area under curve <- 0
 if (length(peaksIni) > 0 && length(peaksFin) > 0) {
   for (j in 1:length(peaksIni)) {
     start <- peaksIni[j]
     end <- ifelse(j <= length(peaksFin), peaksFin[j], NA)</pre>
     if (!is.na(end) && end > start) {
       peak values <- ch4[start:end]
       max values <- sort(peak values, decreasing = TRUE)[1:2]
       avg max values <- mean(max values)
       sum2maxpeaks <- sum2maxpeaks + avg max values
       # Calculate area under the curve above ground level (assumed to be the minimum value in ch4)
       ground level <- min(ch4)
       area under curve <- area under curve + sum(peak values - ground level)
 list(peaksIni = peaksIni, peaksFin = peaksFin, sum2maxpeaks = sum2maxpeaks, area_under_curve = area_under_curve)
! Apply the peak detection to the entire database
UC <- event ch4 %>%
 group by(eventID) %>%
 summarise(
   num peaks = {
     result <- detectPeaksAndSum(ch4)
     length(result$peaksIni)
   sum2maxpeaks = {
     result <- detectPeaksAndSum(ch4)
     result$sum2maxpeaks
   area under curve = {
     result <- detectPeaksAndSum(ch4)
     result$area under curve
```

# Exercise: Calculate all phenotypes for cut visit

```
# Apply the peak detection to the entire database, considering only data between 60 sec to 300 sec per eventID
peak phen2 <- event ch42 %>%
  filter(time diff sec >= 60 & time diff sec <= 300) %>%
  group by(eventID) %>%
  summarise(
   num peaks = {
      result <- detectPeaksAndSum(ch4)
      length(result$peaksIni)
   sum2maxpeaks = {
      result <- detectPeaksAndSum(ch4)
      result$sum2maxpeaks
   area under curve = {
      result <- detectPeaksAndSum(ch4)
      result$area under curve
# Print the results
head(peak phen2)
dim(peak phen2)
```



# Moving average, play with the window size

```
# Function to calculate the average of the moving average of CH4 values
calculate moving avg <- function(ch4, window size) {
  moving avg <- rollmean(ch4, k = window size, fill = NA)
  mean(moving avg, na.rm = TRUE)
# Apply the peak detection to the entire database, considering only data between 60 sec to 300 sec per eventID
all phen <- event ch42 %>%
  filter(time diff sec >= 60 & time diff sec <= 300) %>%
  group by(eventID) %>%
  summarise(
   num peaks = {
      result <- detectPeaksAndSum(ch4)
     length(result$peaksIni)
    sum2maxpeaks = {
      result <- detectPeaksAndSum(ch4)
      result$sum2maxpeaks
    area under curve = {
      result <- detectPeaksAndSum(ch4)
      result$area_under curve
    avg moving avg ch4 = calculate moving avg(ch4, window size = 10) # Change window size as needed
# Print the results
head(all phen)
dim(all phen)
```

#### Exercise: Calc correlations between phenotypes

```
#calculate all the phenotypes with cut and not cut and calculate correlations among the phenotypes
#HINT1: use merge or join
merged<- merge(event_5low, event_quant, by="eventID")
head(merged)
merged2<- left_join(event_5low, event_3low, by="eventID")
head(merged2)
#HINT2: use cor and choose the method kendall, spearman or pearson
cor(event_5lowavg_ch4, event_3lowavg_ch4, method="kendall", use="complete.obs")
cor(event_meansmean_ch4, event_meansmean_co2, method="pearson", use="complete.obs")</pre>
```





# Spain software example



Birgit Gredler-Grandl, Coralia Manzanilla-Pech, Ester Teran and Oscar González-Recio





#### Now let's use an automatic software (Oscar)

- 1. Open settings  $\longrightarrow$  New  $\longrightarrow$  Terminal
- 2. cd DAY1
- load module
   ml use /lustre/shared/Courses/RELIVESTOCK2025/modules
   module load SnifferAnalyzer
- 4. run the program giving the input files or using bash bash.sh sniffer\_analyzer herd\_robot101.csv herd\_sniffer101\_Loggy\_42s\_FD1.txt 42
- 5. check the output.txt





#### Raw Data





Date and time	CH4	CO2
03/12/2021 11:50:00	0,018	0,095
03/12/2021 11:50:01	0,018	0,095
03/12/2021 11:50:02	0,018	0,093
03/12/2021 11:50:03	0,018	0,094
03/12/2021 11:50:04	0,018	0,093
03/12/2021 11:50:05	0,018	0,095
03/12/2021 11:50:06	0,018	0,093
03/12/2021 11:50:07	0,018	0,095
03/12/2021 11:50:08	0,018	0,095
03/12/2021 11:50:09	0,018	0,095
03/12/2021 11:50:10	0,018	0,095
03/12/2021 11:50:11	0,018	0,095
03/12/2021 11:50:12	0,019	0,095
03/12/2021 11:50:13	0,019	0,095
03/12/2021 11:50:14	0,019	0,094
03/12/2021 11:50:15	0,019	0,093
03/12/2021 11:50:16	0,019	0,094
03/12/2021 11:50:17	0,02	0,095
03/12/2021 11:50:18	0,02	0,095
03/12/2021 11:50:19	0,02	0,095
03/12/2021 11:50:20	0,02	0,095
03/12/2021 11:50:21	0,02	0,095
03/12/2021 11:50:22	0,02	0,097
03/12/2021 11:50:23	0,02	0,096
03/12/2021 11:50:24	0,021	0,095
03/12/2021 11:50:25	0,021	0,095

#### Sincronize



SUAREZ\_102

Numero vaca	Addres	Fecha/Hora de visita	tiempo en cubiculo	Produccion de leche	Tiempo	Descripcion	CIB
161	102	22/11/2021 0:02:00	4:26	18.8	6:04	Correcto	ES031112194660
148	102	22/11/2021 0:15:00	10:08	10.9	11:39	Correcto	ES021112642431
193	102	22/11/2021 0:23:00	6:01	15.8	7:38	Correcto	ES091111967021
110	102	22/11/2021 0:28:00	3:45	14.4	5:04	Correcto	ES041112449478
136	102	22/11/2021 0:35:00	5:37	15.1	7:11	Correcto	ES091112194713
158	102	22/11/2021 0:43:00	5:42	15.9	7:29	Correcto	ES031112449535
257	102	22/11/2021 0:48:00	3:41	11.3	5:20	Correcto	ES091112642438
178	102	22/11/2021 0:54:00	3:44	14.6	5:24	Correcto	ES061112449527
142	102	22/11/2021 1:01:00	5:01	15	6:35	Correcto	ES031112194693



									BUAREZ_100ourps	it:						
CH	***	date	high	Tre	пинесия	*********	******CHLC02	peeks	Store of PosterCH4	Bern of Postscoo	State of Posts/Estin	MIK CHI	AUC PARK	Mean of Poskethia	Main, et Pealss Ratio	yeldty
E8041112440029	230	Mon Nev 22 15:48:33 CET 2021	151	398.0	388.82	8224.3	0.08		1482.8071148498303	36504.79056912293	0.200900006440180733	9978.90	1.37	247.1249191408697	0.00468042573084468	.0
E9071112449625	121	Mini Nev 22 19:58 23 CET 2021	0.7	216.0	143.55	4209.27	q.gr	3	382 1966194360681	13815.543683868752	0.14040237529072290	4575.02	1.8	130,7222264790329	0.046467408426674216	. 11
E8021112449545	766	Mon Nov 22:20:09:54 CET 2021	15.7	267.0	77.89	8027.68	0.07	- 8	227.29605654341737	11560 271800728290	0.1257196011929371	2246.25	1.19	75,76645284780576	0.04190955873099909	10
ES041111906965	275	Mon Nov 22 20:13:10 CET 2021	12.0	301.0	15881	3602.45	0.01		485,0052902066624	18271.33913505767	0.5307052387703151	1715.61	4.02	97.00105864138324	0.10014104855409302	10
E9011112449633	368	Mon Nev 20:30:19:49 CET 2001	12.4	322.0	202.1	1000.18	4.00		943.7002000030722	35408.7870015873	0.11710000032222200	6923.54	0.8	168.7452990001446	0.003429016704690794	. 0
E3061112449005	100	Mon Nev 22 2025-15 DET 2021	0.0	410.0	172.68	4002.3	0.09		990.0901341991343	30007 80104805190	0.20150701304008717	5467.40	1.23	166,1160023688204	0.000884038498014530	0
C9061112449635	30	Mon Nov 22 20:34 27 DET 2021	14.2	284.0	58,64	4541.71	0.13	_	29.7752303330754	1414.1570083707866	6.10707299876484006	1410.4E	4.0	39,7752908908764	0.10767290670484006	10
R9081112194724	258	Mon New 22/20/43 66 CET 2021		312.0	129.01	310 0	0.09	П	4 5.83832 72335	19706.5511600094908	0.17946807833184026	4218.08	1.7	113.29868086081349	0.040000119927899276	- 1
608061112A406INB	308	Mon New 22: 20:48 18 CET 2021	20.0	353.0	162.5	364 61	0.07	Р,	5 7,98236 460626	19021.88009974136	03352943296704917	4330,2	2.58	118.88947309012925	0.04006496873409533	10
E5041112194884	125	Mon Nev 20: 20:58:29 CST 2021	160	340.0	94.08	4.78		4	20 10 178501	1000-1400000005900	0.4332133999963586	3949.13	3.33	82.448/0279/012794	p.18890827517583064	. 0
E9081112642404	304	Mon Nov 50 51 55 19 CST 2021	17.7	363.0	124.83	3963.61	0.12	3	416.18420521035106	11718-46149000323	0.16710007950613496	2575.07	1.12	120,720068336764	0.005700000000711654	10
ES071112842447	355	Mon Nov 22 21:33 45 DET 2021	10,4	380.0	143.52	9907.58	41		600.1588758059027	20041.780010728847	0.38961811552911004	6160.1	3.95	107,03177470110663	0.07790922910582205	10
E5051112642396	179	Mon Nev 22 21 42 58 DET 2021	12.1	903.0	197.48	3412.15	9.29	,	1207.601086240045	21018.30443063256	1.350011447947999	6426.87	9.85	172.51458946414928	D.1802879487068527	. 0
E3081112842437	705	Mon Nev 22 21.64.12 CET 2021	12.0	299.0	9529	3022.13	4.11	1	211.11220011229018	9635-074944074843	0.00007870010197078	2180/76	1.00	106,68613805610506	0.00416608797666038	- 0
E3001111107020	232	Min Nev 22 22:00:32 CET 2021	12.4	450.0	68.38	1708.48	0.12		349.90201789023254	11015.28748936095	0.42803751793639104	3806.80	4.5	69,98879477646948	b 1600007000007377221	0
E801111219471S	248	Mon Nov 22 22-18 80 GET 2021	21.5	306.0	202.21	8241,77	0.09		866.96899689687	815151515151,75988	0.2342820583159971	6576.80	2.68	216,74342424242425	D.01007351482994928	10
E90111119000C3	196	Mon Nev 22 32:18:43 CET 2021	11.1	268.0	225.28	5748.11	0.06	4	989.2437160540909	25125.0n6420361348	8:17664566987185848	6877.46	2.00	142.31090901381622	0.043886415417813385	10
E9051112194984	100	Mon Nov 22 32:32:57 CET 0021	12.8	474.D	146.17	6118.57	494	11	1472.3830729111376	72146-20101795497	D 180429000 (08.5346)	6291-66	120	A 7	TW W	
E5091112466542	108	Mon Nev 32 33:54 St OCT 2021	13.2	307.0	120.16	4758-58	4.04	4	176.9674621706736	22/03 72008:000414	0.154509009530031183	1536.20	1.0			/\
E9091112449695	158	Mon Nov 22 22:43:15 DCT 2021	18.2	410.0	100.02	2907.8	0.09	4	434.45247897542264	15196-307130-075158	0.29006/0529160906	4414.04	2.5			A
E509111119666E2	248	Man Nev 22 20:52 50 DET 2021	15,4	871.0	15475	\$155.5	411		E75.4542748194419	E4473.33101922237	0,5456756884758521	4150.94	2.15		TAT	T 3



# Data processing

Java Program	Function
SnifferAnalyzer	<ol> <li>Assign events to cows according to time footprint and lag</li> </ol>
	2. Calculate background (average of 5 lowest measurements from opening of the AMS gate to cow exit).
	3. Detect eructation peaks
	4. Calculate traits
	5. Write output





## Executing SnifferAnalyzer Program

./SnifferAnalyzer\_linux <AMS\_traffic\_file> <sniffer\_file> <lag\_time>



'cnifferAnalyzer\_linux herd\_robot101.csv herd\_sniffer101\_Loggy\_42s\_FD1.txt 42





# What does SnifferAnalyzer do

	A	В		C		D		E		F	G		Н	
1	cow	Address	₩	time ▼	tir	ne 🖪	7	milk 🔻	ti	ime 🔻	Descripcion	-	CIB	W
2	59		101	18/02/2022 0:05:31		05:0	3	15,4	1	05:03	Correcto			59
3	236	1	101	18/02/2022 0:11:10		04:5	7	13	3	04:57	Correcto			236
4	31	1	101	18/02/2022 0:16:33		05:1	8	15,1	L	05:18	Correcto			31
5	231	1	101	18/02/2022 0:26:45		09:2	1	18,7	7	09:21	Correcto			231
6	532	1	101	18/02/2022 0:37:19		05:4	4	14,8	3	05:44	Correcto			532
7	240	1	101	18/02/2022 0:43:43		05:3	1	17,8	3	05:31	Correcto			240
8	259	1	101	18/02/2022 0:52:21		06:4	2	11,8	3	06:42	Correcto			259
9	208		101	18/02/2022 0:58:00		06:0	9	0,7	7	06:09	Tiempo de conexiÔøΩn			208

A	A	В	С
1	Date and time	CH4 ▼	CO2 🔻
325	18/02/2022 00:05:24	0,032	0,534
326	18/02/2022 00:05:25	0,032	0,531
327	18/02/2022 00:05:26	0,032	0,527
328	18/02/2022 00:05:27	0,031	0,514
329	18/02/2022 00:05:28	0,03	0,385
330	18/02/2022 00:05:29	0,03	0,335
331	18/02/2022 00:05:30	0,03	0,331
332	18/02/2022 00:05:31	0,029	0,318
333	18/02/2022 00:05:32	0,029	0,298
334	18/02/2022 00:05:33	0,028	0,296
335	18/02/2022 00:05:34	0,028	0,3
336	18/02/2022 00:05:35	0,029	0,572
337	18/02/2022 00:05:36	0,029	1.002
338	18/02/2022 00:05:37	0,029	1.456
339	18/02/2022 00:05:38	0,03	1.663
340	18/02/2022 00:05:39	0,03	1,61
341	18/02/2022 00:05:40	0,03	1.553
342	18/02/2022 00:05:41	0,03	1,49
343	18/02/2022 00:05:42	0,029	1.438
344	18/02/2022 00:05:43	0,029	1.142
345	18/02/2022 00:05:44	0,029	0,955
346	18/02/2022 00:05:45	0,029	0,717
347	18/02/2022 00:05:46	0,028	0,661
348	18/02/2022 00:05:47	0,028	0,623
349	18/02/2022 00:05:48	0,028	0,597
350	18/02/2022 00:05:49	0,027	0,579
351	18/02/2022 00:05:50	0,027	0,566
352	18/02/2022 00:05:51	0,027	0,564
353	18/02/2022 00:05:52	0,027	0,577
354	18/02/2022 00:05:53	0,027	0,572
355	18/02/2022 00:05:54	0,028	0,567

Instituto Nacional de Investigación y Tecnología Agraria y Alimentaria



#### Data management: Sniffer output



```
Date and time; Input 1 (2) [V/mA] 
2023-11-06 12:29:00; 0,038 
2023-11-06 12:29:01; 0,038 
2023-11-06 12:29:02; 0,038 
2023-11-06 12:29:03; 0,038 
2023-11-06 12:29:04; 0,038 
2023-11-06 12:29:05; 0,039 
2023-11-06 12:29:06; 0,039 
2023-11-06 12:29:07; 0,039 
2023-11-06 12:29:08; 0,039 
2023-11-06 12:29:09; 0,039 
2023-11-06 12:29:09; 0,039 
2023-11-06 12:29:10; 0,039
```





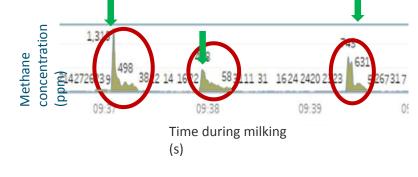




#### METHANE PHENOTYPING

Phenotype definition (weekly averages)

- 1. Mean CH<sub>4</sub> (by second and every 5 s)
- 2. Sum of peaks CH<sub>4</sub> (by second and every 5 s)
- 3. Sum of max peaks CH<sub>4</sub>
- 4. Area under the curve (AUC CH<sub>4</sub>)
- 5. Ratio of (mean) CH<sub>4</sub>/CO<sub>2</sub>
- 6. CH<sub>4</sub> grams per day (Madsen et al., 2010)



Prod CH<sub>4</sub>
$$\left(\frac{g}{d}\right) = 0.714 * ratio(ppm) * 180 * 24 * 0.001 * (5.6 * kg body mass0.75) +22 * ECM + 1.6*10-5 * days in pregnancy$$



7. CH<sub>4</sub> grams per day (in-house\*, based on observed CO2 and CH4/CO2 ratio, to be published)

Large genetic correlation between phenotypes





#### How to clone files/folders from Github to Jupyter

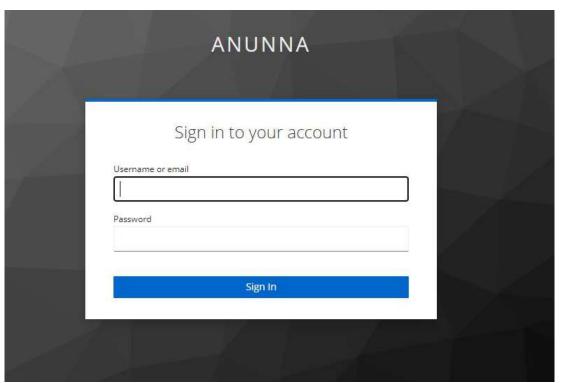
- Open Jupyter Notebook online in your web browser (<u>https://notebook.anunna.wur.nl</u>).
- $2.\,\,$  Click on the "New" button in the top right corner of the screen.
- 3. Select "Terminal" from the dropdown menu.
- 4. In the terminal window, type the following command:

git clone https://github.com/ogrecio/RelivestockMethaneCourse



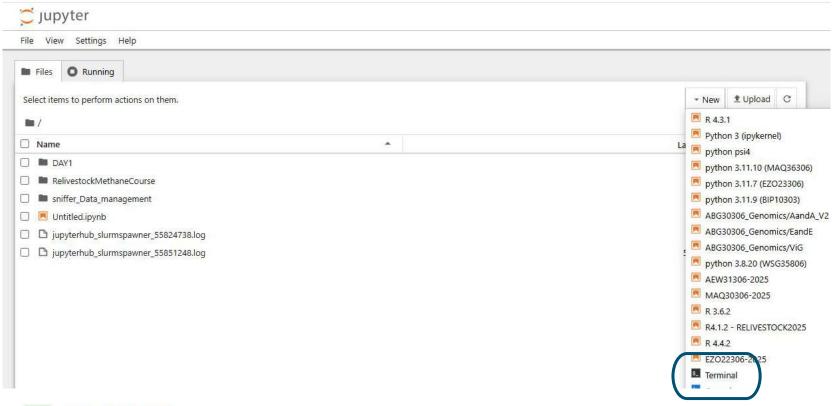


# Step 1- Open Jupyter Notebook and login





# Step 2 and 3- Open the terminal





# Step 4- Clone the repository of the course



File View Settings Help

```
teran01@node217:~$ git clone https://github.com/ogrecio/RelivestockMethaneCourse Cloning into 'RelivestockMethaneCourse'...
remote: Enumerating objects: 275, done.
remote: Counting objects: 100% (137/137), done.
remote: Compressing objects: 100% (113/113), done.
remote: Total 275 (delta 79), reused 23 (delta 23), pack-reused 138 (from 1)
Receiving objects: 100% (275/275), 47.01 MiB | 32.64 MiB/s, done.
Resolving deltas: 100% (120/120), done.
Updating files: 100% (37/37), done.
```



#### Set the working directory

- cd RelivestockMethaneCourse/Day1
- Is

```
Updating files: 100% (37/37), done.

teran01@node217:~$ ls

DAY1 RelivestockMethaneCourse Untitled.ipynb jupyterhub_slurmspawner_55824738.log jupyterhub_slurmspawner_55851248.log sniffer_Data_management teran01@node217:~$ cd RelivestockMethaneCourse $ ls

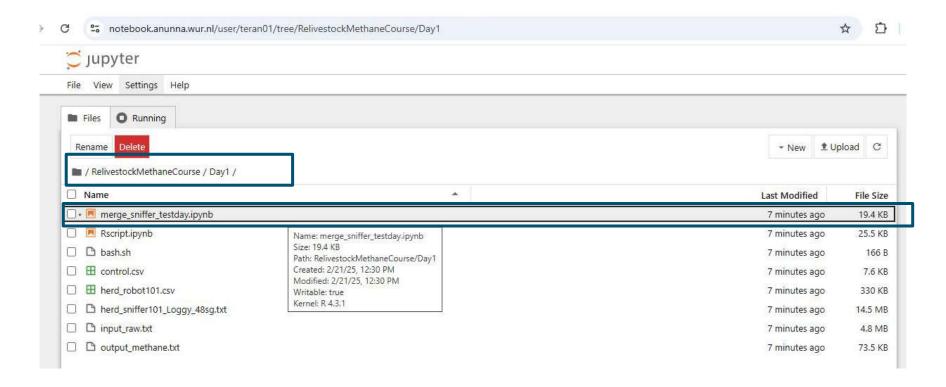
ASRem1-4.2-Functional-Specification.pdf Day1 LICENSE Rscript.ipynb Timetable-ing_v2_Zaragoza_Final.pdf bash.sh programs scripts

ASRem1-4.2-Structural-Specification.pdf Day2 README.md Rscripts __config.yml data reading_rawCH4data.R slides teran01@node217:~{RelivestockMethaneCourse}$ cd Day1 teran01@node217:~{RelivestockMethaneCourse}$ ls

Rscript.ipynb bash.sh control.csv herd_robot101.csv herd_sniffer101_Loggy_48sg.txt input_raw.txt merge_sniffer_testday.ipynb output_methane.txt teran01@node217:~{RelivestockMethaneCourse/Day1}$ |
```



# Back to Notebook and open the script





# Merge\_sniffer\_testday

```
25 notebook.anunna.wur.nl/user/teran01/notebooks/RelivestockMethaneCourse/Day1/merge_sniffer_testday.ipynb
Jupyter merge_sniffer_testday Last Checkpoint: 9 minutes ago
                                                                                                                                                      Not Trusted
File Edit View Run Kernel Settings Help
B + % (1) (1) >
                                                                                                                                     JupyterLab ☐ # R 4.3.1 ○ ■
     [ ]: # Script methane data management
           library(data.table)
           library(dplyr)
           library(plyr)
           library(tidyr)
           library(lubridate)
           library(ggplot2)
           library(GGally)
     [ ]: # Set the working directory where the files are located
           setwd("../data")
```



#### Data management: R

Call the libraries and set the work direction

```
In [ ]:
 # Script methane data management
 library(data.table)
 library(dplyr)
 library(plyr)
 library(tidyr)
 library(lubridate)
 library(ggplot2)
 library(GGally)
```





# Read files and check they were correctly imported

```
In [ ]:

# Read the sniffer output
bd=read.table("output_methane.txt",sep=",",header=T) #500 rows & 20 cols
bd_id=unique(bd[,1]) #63 individuals
summary(bd)
```





## Data management: Output format

- ❖ One farm 1 robot
- 500 records per visit
- **♦** 63 cows

#### **Phenotypes**



-	cow	date	kgm	time	meanCH4	meanCO2	meanRatioCH4_CO2	AUC_CH4	AUC_Ratio	peaks	peaks_per_minu
1	5425	Tue Aug 09 17:31:52 CEST 2022	15.84	327	146.47	2313.80	0.09	14653,21	8.39	5	2
2	6502	Tue Aug 09 17:37:21 CEST 2022	14.39	358	395.57	4347.48	0.09	30028.49	5,72	3	
3	6734	Tue Aug 09 17:43:21 CEST 2022	13.85	418	145.07	2814.31	0,05	5263.64	1,67	3	
4	6493	Tue Aug 09 17:50:21 CEST 2022	13.72	418	235,07	2594.37	0.09	14608.71	4.72	3	
5	7431	Tue Aug 09 17:59:40 CEST 2022	15.86	459	426.97	4366.50	0,08	27266.67	5.05	5	
6	9765	Tue Aug 09 18:09:57 CEST 2022	10.55	322	194.63	2952.69	0.04	19926.71	4.04	4	
7	2504	Tue Aug 09 18:15:26 CEST 2022	14.31	353	245.03	6197.70	0,05	19440.00	3.07	4	
8	2514	Tue Aug 09 18:21:21 CEST 2022	18.76	418	621.93	6548.15	0.11	47065.55	8.07	4	
9	7436	Tue Aug 09 18:34:53 CEST 2022	8.34	266	493.27	6448.67	0,09	48298.65	7,55	4	
10	9750	Tue Aug 09 18:39:21 CEST 2022	18.33	418	375,19	4027.56	0.10	21424.02	6,02	3	
11	7441	Tue Aug 09 18:48:03 CEST 2022	9.56	496	261.73	4234.37	0.04	17648.95	2.51	5	





#### Read files and check they were correctly imported

```
# Read test day file
test=read.table("control.csv",sep=";",header=T)#107 rows
test_id=unique(test[,1])#64 #july and sept
#obtain date from test day records
testtest_date1 = dmy(testtest_date)
# Obtain kgm of protein and fat

testkgmfat = testfat*test$milk/100
testkgmprotein = testprotein*test$milk/100
```





#### Data management: data test description

Test date

To join to sniffer data

Productive traits

Genetic correlations

N of calving

Model

Weight

Methane production

Calving date

Days in milking

*	cow	test_date	numpar	calving_date	milk	fat	protein	RCS	ETS	lactose	urea	bhb	weight	t
1	586	16/09/2022	2	01/08/2022 0:00	39.32	3.34	3.00	106	8.49	4.73	121	0.05	571,0743	1
2	635	19/07/2022	1	NA	21.20	3.93	3.15	74	8.94	5.06	285	0.04	NA	
3	635	16/09/2022	1	NA	21.60	3.92	3.32	40	9.07	5.03	129	NULL	NA	
4	4109	19/07/2022	1	NA	26.80	3.96	3.26	18	8.88	4.95	329	0.04	NA	
5	4109	16/09/2022	1	NA	26.00	3,57	3.53	15	9.24	4.84	149	0.05	NA	
6	4110	19/07/2022	1	07/07/2022 0:00	23.23	3.26	3.12	51	8,56	4.61	200	0.07	598.4168	
7	4110	16/09/2022	1	07/07/2022 0:00	25.74	5,06	3.37	19	9.04	4.99	155	NULL	598.4168	
8	5422	16/09/2022	3	16/08/2022 0:00	43.31	3.66	2.83	14	8,47	4.89	108	0.1	580.8212	
9	5423	16/09/2022	4	NA	30.40	4.67	3.18	11	9.00	5.00	130	0.05	NA	
0	5424	16/09/2022	3	31/01/2022 0:00	28,43	3.84	3.29	173	8.79	4.63	155	0.07	514.0000	
4	5/12/	10/07/2022	2	31/01/2022 0:00	37.66	2.74	202	2/1	8.61	4 00	265	0.04	514 0000	





# Join sniffer file with test file by ID and closest date

```
# Obtain date from output
bd<-tidyr::separate(data = bd, col ="date",into = c("day", "month", "date", "time", "z", "year"),
                                                                                                       Check the format
                   sep = " ", extra = "merge")
                                                                                                       of date column
bd$sniffer date=paste(bd$date, bd$month, bd$year, sep = "/")
bd$sniffer date1<-dmy(bd$sniffer date)
# Join the two datatables-
bd full <- lapply(intersect(bd$cow,test$cow),function(id) {
 d1 <- subset(bd,cow==id)
 d2 <- subset(test,cow==id)
 d1$indices <- sapply(as.Date(d1$sniffer date1),function(d) which.min(abs(as.Date(d2$test date1) - d)))
 d2$indices <- 1:nrow(d2)
 merge(d1,d2,by=c('cow', 'indices'))
}) ###
bd full <- do.call(rbind,bd full) #342
bd full$indices <- NULL
```





# Check the number of animals that have not joined





#### Days in milking and threshold

```
# Obtain days in milking
bd_full=tidyr::separate(data=bd_full,col ="calving_date", into=c('calving_date', 'calving_time'),sep=' ')
bd_full$calving_date1=dmy(bd_full$calving_date)
bd_full$daysinmilking=bd_full$sniffer_date1-bd_full$calving_date1
bd_full$daysinmilking=gsub("[a-z]","",bd_full$daysinmilking)#retain only numbers
bd_full$daysinmilking=as.numeric(bd_full$daysinmilking)
# Check difference of days
table(bd_full$daysinmilking)
# Choose The threshold the difference of days , 365
bd_full1=bd_full %>% filter(daysinmilking <= 365)</pre>
```





#### Week of lactation and/or state of lactation

```
# Obtain week of lactation
bd_full1$week_lactation=floor(bd_full1$daysinmilking / 7)
# Codify days in milking as levels
bd_full1<-mutate(bd_full1,state_lactation=case_when(
    daysinmilking < 91 ~ "1",
    daysinmilking > 90 & daysinmilking < 151 ~ "2",
    daysinmilking >150 ~ "3"))
table(bd_full1$state_lactation)
```





#### Number of calving

```
# Codify number of calving in three or four levels
hist(bd_full1$numpar)
table(bd_full1$numpar)
# If there is data of more than 3 calvings
bd_full1<-mutate(bd_full1,num_calving=case_when(
    numpar <= 1 ~ "1",
    numpar > 1 & numpar < 3 ~ "2",
    numpar > 2 ~ "3"))
table(bd_full1$num_calving)
```



# Data management: data output description

Methane phenotypes : ppm



- Ratio CH4/CO2
- Methane production: grams per day (Madsen et.al, 2010)

  Prod CH<sub>4</sub>(g/d)=0.714\*ratio\*180\*24\*0.001\*(5.6\*kg body mass<sup>0.75)</sup>+

  22\*ECM +1.6\*10<sup>-5</sup>\*days in pregnancy
- Methane production: Tier 2. Dependent of milk, fat, protein and weight



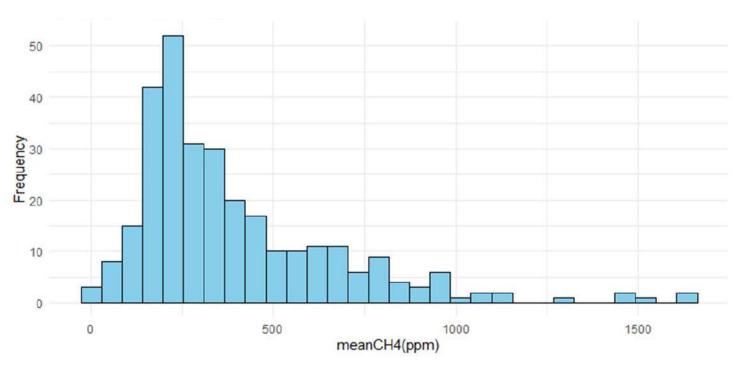


#### Ratio and grams per day

```
# Ratio mean CH4/CO2
bd full1$ratioCH4CO2=bd full1$meanCH4/bd full1$meanCO2
mean(bd full1$ratioCH4CO2)
# Obtain grams per day
# Madsen equation
bd full1$ECM=bd full1$milk*(0.25+0.122*bd full1$fat+0.077*bd full1$protein)
bd full1$days inpregnancy=0
bd full1$gd madsen<-(0.714*bd full1$ratioCH4CO2)*180*24*0.001*(5.6*bd full1$weight**0.75+22*bd full1$ECM+
                                                                 1.6*0.00001*bd full1$days inpregnancy**3)
mean(bd full1$gd madsen,na.rm=T)
# Tier 2 equation
Cf<-0.386*1.2 #coefficient for feeding situation, lactating
Ca<-0 #Coefficient for activity
Cp<-0.10 #coefficient for pregnancy
DE<-80 #Digestible energy per gross energy in cows fed low quality forage and concentrate
NEm<-Cf*bd full1$weight**0.75
NEak-NEm*Ca
NEl<-bd full1$milk*(1.47+0.4*bd full1$fat)
NEp<-ifelse(bd full1$daysinmilking>100, Cp*NEm,0)
REM<-1.123-(0.004092)*DE+0.00001126*(DE)**2-25.4/DE
GE<- 1000*((NEm+NEa+NEl+NEp)/REM)/(DE/100) #g/d
bd full1$CH4 tier2<-GE*0.065/55.65
mean(bd full1$CH4 tier2, na.rm=T)
```



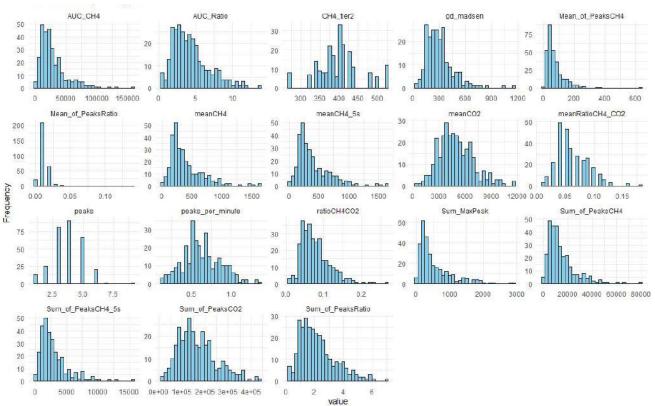
## Data management: Phenotypes distribution







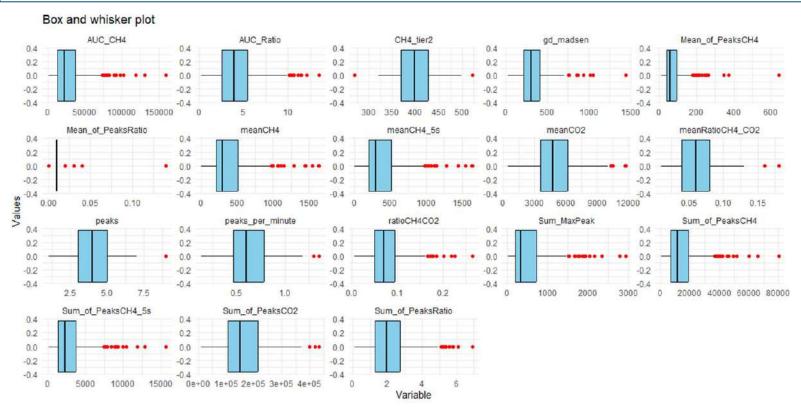
#### Data management: Phenotypes distribution







#### Data management: Outlier visualization







# Data management: filtering

- CO2 concentration lower limit, as technical error < 2500
- Filter applied to ratio CH4/CO2 > 0.3
- Filter applied to CH4 g/d > 1200 values
- Detection of outliers
- Data correction of outliers ± 3SD





#### Filtering out

- Methane records CO2<2500.</li>
  - > +- 3SD in the comparison group
  - <5 or >35 records per week (remove all records in that week)

Remove weeks with <5 records, and animals with only 1 week of records

- CO2 records
  - <2500
  - > +- 3SD in the comparison group
- CH4/CO2 records ratio>0.3
  - > +- 3SD in the comparison group

Average data per week (day), for consistency (more reliable heritability estimates)





#### Data management: filtering

Set as NAs the values in the different phenotypes where the corresponding meanCO2 is lower than 2500





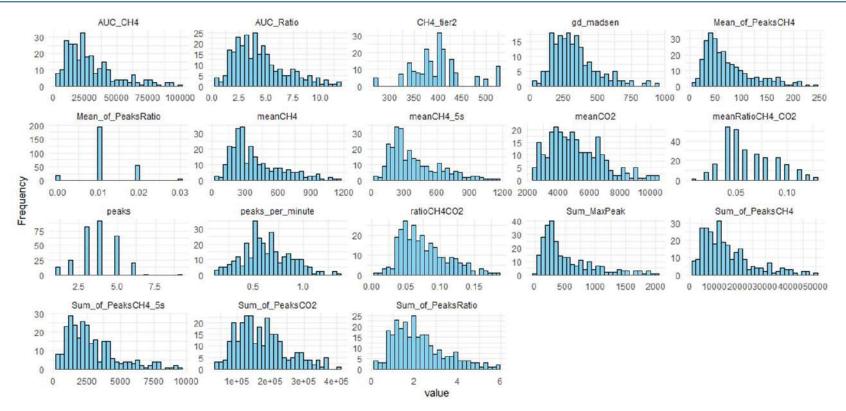
# Data management: filtering

```
# Function for correction of data ±3 SD
function outliers <- function(x) {
 mean x \leftarrow mean(x, na.rm = TRUE)
 sd \times (-sd(x, na.rm = TRUE))
 # Set upper and lower lim (± 3 SD)
 lower lim <- mean x - 3 * sd x
 upper \lim x - \max x + 3 * sd x
 # Replace out-of-limit values with Nas
 x <- ifelse(x < lower lim | x > upper lim, NA, x)
 return(x)
#Apply the function
data corrected <- bd full2 %>%
 dplyr::mutate(across(
    c(meanCH4, meanCH4 5s, meanCO2, meanRatioCH4 CO2, AUC CH4, AUC Ratio,
      Sum of PeaksCH4, Sum of PeaksCH4 5s, Sum of PeaksCO2, Sum of PeaksRatio,
      Mean of PeaksCH4, Mean of PeaksRatio, Sum MaxPeak, ratioCH4CO2, gd madsen,
      CH4 tier2, milk, kgmfat, kgmprotein), #Apply the function to these variables
    function outliers # Function to replace outliers with NAs
colSums(is.na(data corrected))
```





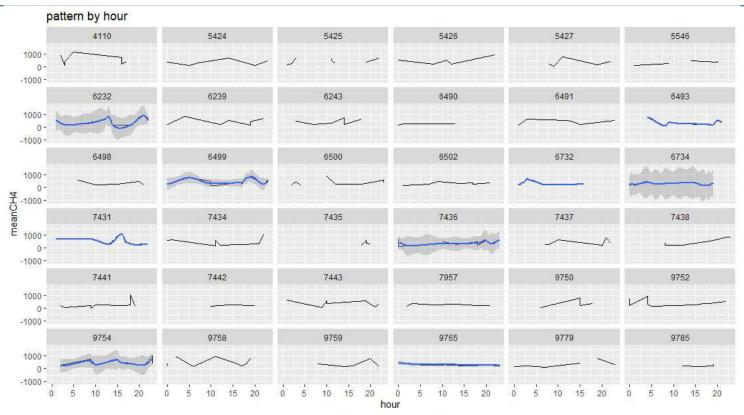
## Data management: data distribution after filtering







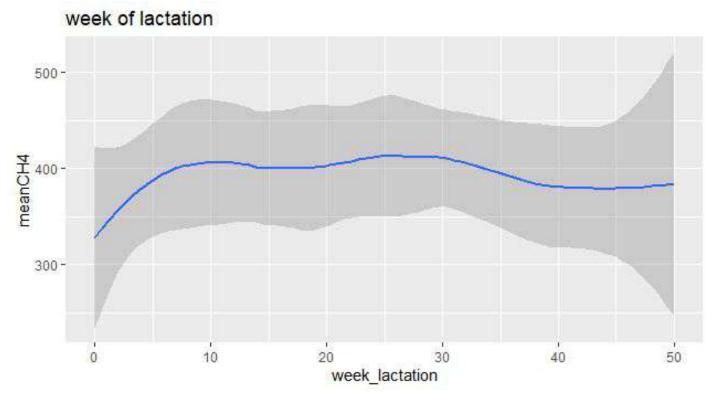
# Data management:pattern by hour







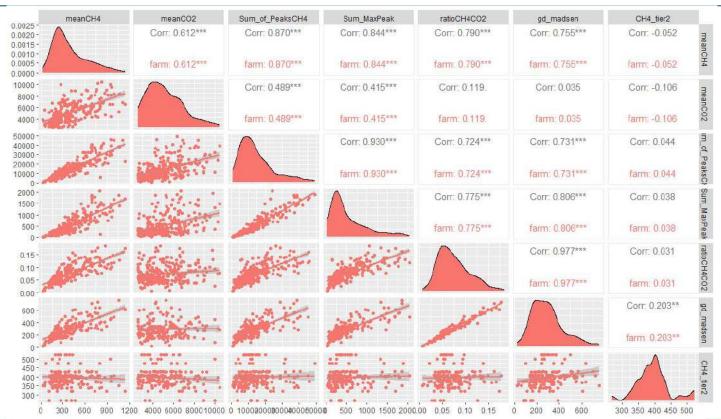
#### Data management: trend by week of lactation







#### Correlation between phenotypes







- Data more consistent.
- Check the number of data by individual and comparison group (At least 2 per individual and 3/5 per group).



- Data more consistent.
- Check the number of data by individual and comparison group (At least 2 per individual and 3/5 per group).

```
# Calculate the mean per week per animal, considering robot and year if there is more than one mean_week <- data corrected %>%

group_b/(cow,epiweek,epiyear) %>% #some farms could be studied in two or more periods dplyr::summarise(across((where(is.numeric)), list(mean = ~mean(.,na.rm=T)), .names = "{.col}_r count = n()) %>%

dplyr::left_join(data_corrected, by = c("cow","epiweek","epiyear"))#keep the other variables (are necessary fo # Filter by count per week, keep count > 3/5 records pe week table(mean_week$count)

mean_week1=mean_week[mean_week$count > 3, ]
```

Keep the information for further analysis





```
# Choose the variables to keep
names(mean_week1)
mean_week2=mean_week1[,c(1:3,5:38,86,87,94,95)] # ID, averages of phenotypes and variables to further analysis
mean_week3=mean_week2[!duplicated(mean_week2),]
# Sort the columns
names(mean_week3)
mean_week4=mean_week3[,c(1:3,41,40,39,38,4:30,33,36,37)] #38 rows
```



```
# Check the number of animals by comparison group. i.e =herd-season-year or herd-robot-week-year

vmean_week5=mean_week5%>%
    group_by(farm,robot,epiweek,epiyear)%>%
    dplyr::mutate(N_group=n())

table(mean_week5$N_group)

# Check the number of records per animal (to be considered in the model, repeated measurements)

vmean_week5=mean_week5%>%
    group_by(cow)%>%
    dplyr::mutate(N_group=n())

table(mean_week5$N_anim)
```





# Data management: Final

^	cow	epiweek *	epiyear	state_lactation	num_calving	robot	farm <sup>©</sup>	meanCH4_mean_week	meanCO2_mean_week	meanRatioCH4_CO2
1	4110	32	2022	1	1	1	1	564,4129	6502,296	Δ.
2	5424	32	2022	3	3	1	1	348.6350	4372.060	
3	5425	32	2022	3	3	1	1	444.1787	4900.474	
4	5426	32	2022	2	3	1	1	493,4380	5193.396	
5	5427	32	2022	3	3	1	1	321.1300	4112,284	
6	5546	32	2022	1	1	1	1	321.0280	6483,494	
7	6232	32	2022	2	2	1	1	481.0820	5271,495	
8	6239	32	2022	1	3	1	1	435.8071	6298,911	
9	6243	32	2022	3	3	1	1	484.4729	6602,339	
10	6490	32	2022	2	2	1	1	398.3675	5782.205	







**Effects** 

Phenotypes

