

Statistical Analysis of Eye Tracking Data

Krzysztof Krejtz

LEAD-ME Summer Training School Warsaw 2021 (5-9 July 2021)  
Eye tracking in media accessibility research - methods, technologies and data analyses

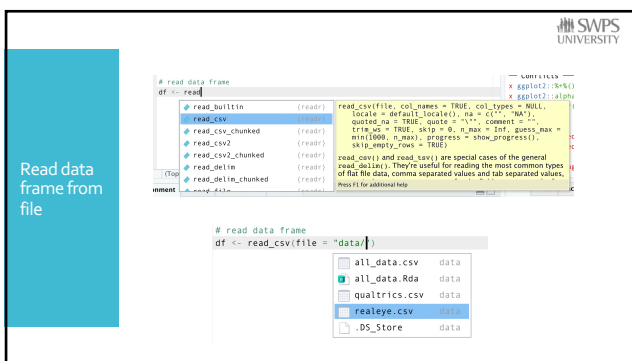
1



Prepare data for the analyses

read / filter / select / mutate / join / write

2



Read data frame from file

```
# read data frame
df <- read_csv("data/real_eyedata.csv")
```

Column specification

ID = col\_double(),  
aot\_id = col\_character(),  
aot\_name = col\_character(),  
aot\_size\_percents = col\_double(),  
aot\_fixation\_total\_count = col\_double(),  
aot\_fixation\_average\_duration\_ms = col\_double(),  
aot\_fixation\_ttf\_ms = col\_double(),  
aot\_fixation\_average\_total\_time\_spent\_ms = col\_double(),  
aot\_fixation\_first\_fixation\_average\_duration\_ms = col\_double(),  
notes = col\_character()

3



Read data frame from file

```
> df <- read_csv(file = "data/real_eyedata.csv")
```

Column specification

ID = col\_double(),  
aot\_id = col\_character(),  
aot\_name = col\_character(),  
aot\_size\_percents = col\_double(),  
aot\_fixation\_total\_count = col\_double(),  
aot\_fixation\_average\_duration\_ms = col\_double(),  
aot\_fixation\_ttf\_ms = col\_double(),  
aot\_fixation\_average\_total\_time\_spent\_ms = col\_double(),  
aot\_fixation\_first\_fixation\_average\_duration\_ms = col\_double(),  
notes = col\_character()

4

First look at data

5

Create AOI type independent variable

- There were three different types of AOI (positive picture, negative and neutral)
- Currently data frame contains AOI names as follows:

```
> unique(df$aoi_name)
[1] "negatyw_pozar"      "neutral_igla"      "pozytyw_ocean"     "negatyw_lodowiec"
[5] "neutral_kosci"      "pozytyw_gory"      "negatyw_niedzwiedz" "neutral_lyzeczeki"
[9] "pozytyw_las"        "negatyw_podtopienia" "neutral_muszl"      "pozytyw_morze"
[13] "negatyw_susza"      "neutral_zegar"      "pozytyw_rzeka"
```

- We want to make AOI type variable which will be used as a factor in further analysis
- We will use first element of unique AOI names as values of this new variable
- Technically we will split aoi\_name variable into two columns by "\_", creating new "aoi\_type" and "picture" variables

```
df <- df %>%
  separate(aoi_name, sep = "_", into = c("aoi_type", "picture"))
```

6

Select & filter variables

- Sometimes you want to delete from your data frame some variables (columns)
- This time we want to get rid of two variables "notes", "aoi\_id", and "picture"
- Also we want to remove four subjects which were not following the instruction

```
df <- df %>%
  select(-notes, -aoi_id, -picture) %>%
  filter(ID != "61567" & ID != "64828" & ID != "66001" & ID != "54127")
```

- The list of new df variables names shows that the operation was successful

```
> names(df)
[1] "ID"
[3] "aoi_size_percent"
[5] "aoi_fixation_average_duration_ms"
[7] "aoi_fixation_average_total_time_spent_ms"
[9] "aoi_type"
[11] "aoi_fixation_total_count"
[13] "aoi_fixation_ttfb_ms"
[15] "aoi_fixation_first_fixation_average_duration_ms"
```

7

Missing values (NAs) a special case

- In R typically missing values are annotated with NA
- In statistical analysis NAs are not welcome
- Always try to find out why you have NAs in your data set
- This time NAs meaning is that the persons did not fixate an AOI
  - It make sense to replace NAs with 0

```
df <- df %>%
  replace_na(list(aoi_fixation_ttfb_ms = 0, aoi_fixation_average_duration_ms = 0,
                 aoi_fixation_average_total_time_spent_ms = 0,
                 aoi_fixation_first_fixation_average_duration_ms = 0))
```

8

Read new data file with questionnaire answers

```
# read data file with additional variables
dq <- read_csv(file = "data/qualtrics.csv")
```

Column specification

```
cols(
  Q1 = col_double(),
  Q3 = col_double(),
  Q5 = col_double(),
  Q7 = col_double(),
  Q9 = col_double(),
  Q11 = col_double(),
  Q13 = col_double(),
  Q15 = col_double(),
  ID = col_double()
)
```

	Q1	Q3	Q5	Q7	Q9	Q11	Q13	Q15	ID
1	4	5	4	5	5	5	4	4	43428
2	4	5	2	5	4	5	3	2	64666
3	4	4	4	5	4	4	2	4	32185
4	2	2	4	2	5	4	3	3	65126
5	4	4	4	5	3	5	5	5	45652
6	1	5	5	4	5	4	2	5	60466
7	3	5	4	5	4	4	4	4	63348
8	3	4	4	4	4	4	5	4	45940
9	2	4	5	5	2	4	5	5	14039
10	3	4	4	2	4	3	4	3	29236
11	4	2	3	3	5	3	4	4	54619
12	4	4	3	4	2	5	5	4	46437
13	4	4	4	4	5	4	2	4	65905

9

Joining two data frames by subject ID

- We need to join two data frames (with eye tracking data and questionnaire data)
- We will join to data frames by subject ID variable

```
# join two data frames by ID (subject ID variable)
d <- inner_join(df, dq, by = "ID")
```

Environment History Connections Tutorial

df 2484 obs. of 18 variables

df 6826 obs. of 8 variables

dq 48 obs. of 9 variables

- Note that new data frame "d" has lower number of rows.
- It is due to the fact that not all subjects who participated in eye tracking study completed also the questionnaire.

10

Calculate independent variables / factors

- Questionnaire data contains answers from the New Ecological Paradigm (NEP) questionnaire and **subject ID**
  - The answers were on 1-5 Likert-type scale (the higher value the higher sensitivity to climate change)
- We want to calculate one score of the sensitivity to climate change (NEP) which will be a mean of all given answers
- Next, we want to perform median split on NEP score (low vs high sensitivity to climate change) to use it as a factor in further analysis.
- We need also to set aoi\_type as factor and make "neutral" value of it as a reference point (important for statistical analysis)
- Last, we do not need all row answers to each NEP question (they all starts with "Q").

```
d <- d %>%
  mutate(NEP = rowSums(select(., starts_with("Q")))) %>%
  mutate(NEPsplit = cut(NEP, 2, labels = c("low", "high"))) %>%
  mutate(aoi_type = factor(aoi_type)) %>%
  mutate(aoi_type = relevel(aoi_type, ref = "neutral")) %>%
  select(-starts_with("Q"))
```

11

Write data files

- Save the entire data frame into a new file
- We can do it in several formats. The most useful are:
  - .csv files (great for sharing even with those who do not use R)
  - Data format (.Rda) - great for further use within R and hard drive space saver

```
# write the entire data base into RData file
save(d, file = "data/all_data.Rda")
# or csv file
write_csv(d, file = "data/all_data.csv")
```

Files Plots Packages Help Viewer

data

qualtrics.csv 1.1 KB Jul 3, 2021, 8:43 PM

realeye.csv 548 KB Jul 3, 2021, 9:18 PM

all\_data.Rda 24.5 KB Jul 3, 2021, 9:43 PM

all\_data.csv 138.5 KB Jul 3, 2021, 9:44 PM

12

## Descriptive statistics

With visualisations

13

### Create new R script "analysis.R"

- It is good to have several R scripts for different purposes and ease-of-reading.
- We will create new R script named "analysis.R"
- Start the script with useful libraries

```
analysis.R
1 # load useful packages
2 library(psych)
3 library(tidyverse)
4 library(afix)
5 library(emmeans)
6
7
8
9
```

14

### Read data frame of RData format

- Reading RData format will load automatically the data frame name

```
# read data frame from RData format
load("data/all_data.Rda")
```

Environment History Connections Tutorial

Global Environment

Data

d 2438 obs. of 10 variables

	ID	aot_type	aot_size_percent	aot_fixation_total_count	aot_fixation_average_duration_ms	aot_fixation_tfff_ms
1	63148	negative	13.51376	11	135	529
2	63148	neutral	13.51376	0	NA	NA
3	63148	positive	13.46663	2	149	7384
4	45040	negative	13.51376	13	190	1106
5	45040	neutral	13.51376	3	200	747
6	45040	positive	13.46663	12	175	4859
7	14170	negative	13.51376	10	145	707
8	14170	neutral	13.51376	7	137	140
9	14170	positive	13.46663	9	131	1343
10	61167	negative	13.51376	2	137	2292
11	61167	neutral	13.51376	4	148	872

15

### Describe data frame

- Summary is a generic function used to produce result summaries of the results of various model fitting functions.
- The function invokes particular methods which depend on the class of the first argument.
- When applied to data frame it returns basic descriptive statistics for all numerical variables

```
summary(d)
      ID          aot_type      aot_size_percent aot_fixation_total_count aot_fixation_average_duration_ms aot_fixation_tfff_ms
Min.   :18564 Length:2438      Min.   :11.97      Min.   : 0.000      Min.   :180.0      Min.   : 500
1st Qu.:41719 Class:character  1st Qu.:12.18      1st Qu.: 1.000      1st Qu.:104.0      1st Qu.:1044
Median :55152 Mode :character  Median :12.38      Median : 4.000      Median :162.0      Median :1918
Mean   :51814              Mean   :12.46      Mean   : 6.779      Mean   :169.1      Mean   :3119
3rd Qu.:63148              3rd Qu.:12.88      3rd Qu.:18.000      3rd Qu.:183.0      3rd Qu.:3998
Max.   :67201              Max.   :13.52      Max.   :44.000      Max.   :583.0      Max.   :14862
NA's   :486              NA's   :486              NA's   :486              NA's   :486

aot_fixation_average_total_time_spent_ms aot_fixation_first_fixation_average_duration_ms NEP NEPsplit
Min.   : 1000              Min.   : 500.0              Min.   :3.125      low :1215
1st Qu.: 450              1st Qu.: 123.0              1st Qu.:3.750      high:1385
Median :1048              Median : 143.0              Median :4.125
Mean   :1492              Mean   : 163.9              Mean   :4.113
3rd Qu.:2020              3rd Qu.: 179.0              3rd Qu.:4.580
Max.   :13419             Max.   :2851.0              Max.   :5.880
NA's   :486              NA's   :486
```

Can be hard to read however

16

Descriptive statistics of selected variables

The **describe** function in the **psych** package is meant to produce the most frequently requested stats in psychometric and psychology studies, and to produce them in an easy to read data.frame.

```
> describe(A.B1)
```

vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
aol_fixation_total_count	1 2430	6.78	7.42	0.0	5.50	5.33	0	44	44	2.52	2.34	0.15
aol_fixation_average_duration_ms	2 1944	169.89	42.81	162.0	163.86	78.17	100	583	483	3.16	18.95	0.97
aol_fixation_ttf_ms	3 1944	3119.15	3863.71	1339.5	2495.53	1830.88	500	14882	14382	1.80	2.88	69.49
aol_fixation_average_total_time_spent_ms	4 1944	1492.11	1536.47	1048.5	1238.19	1807.44	100	13419	13319	2.86	14.39	34.85
aol_fixation_first_fixation_average_duration_ms	5 1944	163.92	79.41	143.0	151.01	35.18	100	2051	1951	8.71	169.99	1.80

The **fast=TRUE** option will lead to a speed up of about 50% for larger data sets

It is also easy to read!

```
> describe(d1,4:8[, fast = TRUE])
```

vars	n	mean	sd	min	max	range	se
aol_fixation_total_count	1 2430	6.78	7.42	0	44	44	0.15
aol_fixation_average_duration_ms	2 1944	169.89	42.81	100	583	483	0.97
aol_fixation_ttf_ms	3 1944	3119.15	3863.71	500	14882	14382	69.49
aol_fixation_average_total_time_spent_ms	4 1944	1492.11	1536.47	100	13419	13319	34.85
aol_fixation_first_fixation_average_duration_ms	5 1944	163.92	79.41	100	2051	1951	1.80

17

# plot distribution of variable values with histogram

```
hist(d1$aol_fixation_average_duration_ms)
```

# plot distribution of variable values with boxplot

```
boxplot(d1$aol_fixation_average_duration_ms)
```

Useful to visually inspect normality of distribution

Useful to identify outlying values

18

Saving graphs in RStudio graphical interface

19

## Moderation analysis with linear regression

Mixed-design linear regression with continuous and nominal predictors and interaction term

20

Linear regression

- In regression we are fitting a model to our data
  - To not loose variance if predictor is continuous and does not need to be spliced into categorical variable
  - To predict values of the dependent variable from one or more independent variable
  - To understand the relationship
    - How does one variable change as the other changes?
    - How much ...does ..... rise with a one unit increase in .....

21

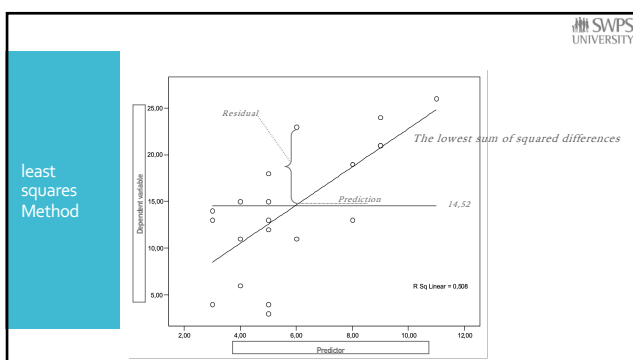
Linear regression

- The outcome variable is predicted using the equation of a straight line for which the squared differences between line and the actual data is minimised.

$$\hat{Y} = a + bX$$

- $\hat{Y}$  = the predicted value of  $Y$  (.....)
- $X$  = ..... (predictor, independent variable)
- Indices, coefficients are computed to asses how accurately  $Y$  scores are predicted by the linear equation:
  - $b$  (slope) - the amount of change in predicted  $Y$  for one unit change in  $X$
  - $a$  (constant) - an intercept, value of  $\hat{Y}$  when  $X = 0$

22




23

Hypotheses testing

- Null hypotheses
  - $b$  slope = 0
  - $a$  constant = 0


24



Simple vs.  
Multiple  
regression

- In simple regression
  - The outcome variable  $Y$  is predicted using the equation of a straight line
- Multiple regression
  - A logical extension of the principles to situation in which there are several predictors.
  - Outcome = (Model) + error


25



Multiple  
regression  
without  
and with  
interaction

- Each predictor has its own regression coefficient
  - For every extra predictor you include, another coefficient need to be estimated
- We are looking for linear combination of predictors that correlate maximally with the outcome variable
  - multiple regression formula with two predictors (no interaction term)
 
$$\hat{Y} = a + \beta_1 X_1 + \beta_2 X_2 + \epsilon$$
  - multiple regression formula with two predictors and interaction term
 
$$\hat{Y} = a + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$$

26




Usefull  
libraries

```
library(emmeans)
library(lme4)
```

**lme4** provides functions for fitting and analyzing mixed models:  
linear (lmer),  
generalized linear (glmer), and  
nonlinear (nlmer.)

27



Hypotheses

- We want to test the following hypotheses:
  - The more sensitive to climate change people are the longer fixation duration on environment pictures
  - Positive and negative pictures of environment will evoke longer fixation duration
  - Sensitivity to climate change will predict fixation duration differently while looking on environment pictures of different emotional valence (interaction hypothesis).

28

**Multiple regression. Model definition**

```
fit <- lmer(aoi_fixation_average_duration_ms ~ NEP + aoi_type + NEP:aoi_type + (1 | ID),
  data = d)
```

Labels in the code above:

- Dependent variable:** `aoi_fixation_average_duration_ms`
- Predictor 1:** `NEP`
- Predictor 2:** `aoi_type`
- Interaction of predictors:** `NEP:aoi_type`
- Random intercept for subjects:** `(1 | ID)`

**Shorter form of full model definition**

```
fit <- lmer(aoi_fixation_average_duration_ms ~ NEP*aoi_type + (1 | ID),
  data = d)
```

29

**Model results basic table**

```
> anova(fit)
```

Type III	Sum Sq	Mean Sq	NumDF	DenDF	F value	Pr(>F)
NEP	2269	2269.4	1	42	0.3989	0.53107
aoi_type	35625	17812.5	2	1932	3.1312	0.04389 *
NEP:aoi_type	51223	25611.7	2	1932	4.5022	0.01120 *

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

30

**Model results. Detailed table**

```
> summary(fit)
```

```
Linear mixed model fit by REML. t-tests use Satterthwaite's method [lmerModLmerTest]
Formula: aoi_fixation_average_duration_ms ~ NEP + aoi_type + NEP:aoi_type + (1 | ID)
Data: d
```

REML criterion at convergence: 22769

Scaled residuals:

Min	1Q	Median	3Q	Max
-2.2882	-0.1972	0.2193	0.5200	1.7394

Random effects:

Groups	Name	Variance	Std. Dev.
ID	(Intercept)	387.1	17.32

Residuals: 1688.8 75.42

Number of obs: 1980, groups: ID, 44

Fixed effects:

	Estimate	Std. Error	df	t value	Pr(> t )
(Intercept)	165.8048	31.1453	184.1182	5.187	<.0005 ***
NEP	-1.0629	0.9374	184.1182	-1.134	0.25946
aoi_typeneutral	-81.3758	32.7634	1932.0000	-2.483	0.01312 *
aoi_typeposityw	-31.8678	32.7634	1932.0000	-0.973	0.33085
NEP:aoi_typeneutral	2.5983	0.9835	1932.0000	2.394	0.02739 **
NEP:aoi_typeposityw	1.6588	0.9835	1932.0000	1.674	0.09424 .

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Correlation of Fixed Effects:

	(Inter) NEP	a_typeposityw	a_typeneutral	NEP:a_typeposityw	NEP:a_typeneutral
NEP	0.993				
a_typeposityw	-0.526	0.521			
a_typeneutral	-0.526	0.521	0.680		
NEP:a_typeposityw	0.521	-0.526	-0.992	0.496	
NEP:a_typeneutral	0.521	-0.526	-0.496	-0.992	0.580

**NOTE:** In the report remember to provide all the values from the coefficients table: coefficient value, standard error, t-test value, degrees of freedom and p-value.

31

**Trends analysis for interaction**

```
em <- emtrends(fit, ~ NEP:aoi_type, var = "NEP")
```

**Code**

```
> print(em)
```

NEP	aoi_type	NEP trend	SE	df	lower CL	upper CL
33	neutral	-1.063	0.937	185	-2.9218	0.796
33	negatyw	1.887	0.937	185	0.0286	3.746
33	pozytyw	0.587	0.937	185	-1.2718	2.446

Degrees-of-freedom method: kenward-roger  
Confidence level used: 0.95

**Results of trends analysis**

**Statistical comparison of slopes**

```
> pairs(em)
```

contrast	estimate	SE	df	t.ratio	p.value
32.977272727272727 neutral - 32.977272727272727 negatyw	-2.95	0.986	1932	-2.994	0.0079
32.977272727272727 neutral - 32.977272727272727 pozytyw	-1.65	0.986	1932	-1.674	0.2154
32.977272727272727 negatyw - 32.977272727272727 pozytyw	1.30	0.986	1932	1.319	0.3845

Degrees-of-freedom method: kenward-roger  
P value adjustment: tukey method for comparing a family of 3 estimates

32

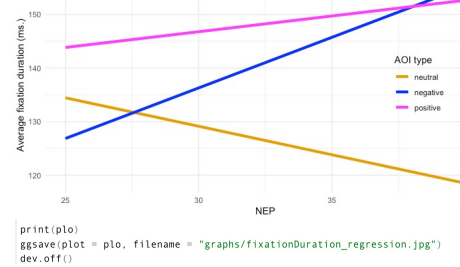
Interaction  
visualization.

Publication  
ready

```
plo <- emmip(fit, aoi_type ~ NEP, cov.reduce = range) +
  geom_line(size=1.5) +
  scale_y_continuous(name = "Average fixation duration (ms.)") +
  scale_x_continuous(name = "NEP") +
  scale_color_manual(name = "AOI type",
    labels = c("neutral", "negative", "positive"),
    values = c("#E69F00", "#0000FF", "#FF00FF")) +
  theme_minimal() +
  theme(legend.position = c(.85, .5))
```

33

Save the  
graph



34

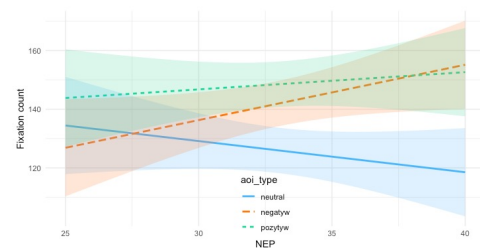
Alternative  
plot of  
interaction

```
# alternatively
require(interactions)
interact_plot(fit, pred = NEP, modx = aoi_type, interval = TRUE) +
  scale_y_continuous(name = "Fixation count") +
  scale_x_continuous(name = "NEP") +
  theme_minimal() +
  theme(legend.position = c(.5, .1))
```

35

Alternative  
plot of  
interaction


(with confidence  
intervals for slopes  
and different  
predefined line  
types)



36

After this workshop ...  
you are able to perform mixed-designed tests of hypotheses for main effects and interaction.  
  
(you should be able to run tests of moderation)


- **Prepare the data frames for the analysis**
  - Read/load data into R
  - Calculate new variables
  - Prepare dichotomous factor with median-split
  - Select variables and filter observations
  - Merge two data frames
  - Write data frame to a file
- **Perform mixed-design ANOVA**
  - Read and interpret the results
  - Estimate means
  - Run pairwise comparisons (post hoc tests) for significant effects
  - Prepare publication ready bar graphs
- **Perform mixed-design linear multiple regression analysis**
  - Read and interpret the results
  - Perform trends (simple slopes) analysis
  - Compare statistically simple slopes
  - Prepare publication ready line graphs with confidence intervals




37

Thank you!

Any questions?





38