

B.L. Experiment Result Analysis

Contents

1	Overview	2
1.1	Summaries of the columns	2
1.2	Reaction time and Mistake Overview	2
2	Mistake Analysis w.r.t. Participant	5
3	Mistake Analysis w.r.t. Stimulus	8
3.1	Stimuli	9
3.2	Image	10
3.3	Category Name	11
3.4	Branch	12
3.5	Level	13
3.6	Stimulus Type	14
4	Component significance analysis in terms of error rates	15
4.1	One-Way	15
4.2	Two-way	17
5	Analysis of reaction times	21
5.1	Overview	21
5.2	Analysis	22
6	B.L. Determination	25
6.1	Simple rating	25
6.2	Probabilistic rating	26
6.3	Probabilistic rating with error	26
7	References	27

1 Overview

1.1 Summaries of the columns

```
##                               stimulus                               lvl
## clothing_n_01|dresser_n_05      : 50  hypernym:3300
## hand_tool_n_01|highboy_n_01     : 50  bl      :3300
## musical_instrument_n_01|armchair_n_01 : 50  hyponym :3300
## acoustic_guitar_n_01|acoustic_guitar_n_01: 25
## acoustic_guitar_n_01|damson_n_01   : 25
## apple_n_01|granny_smith_n_01      : 25
## (Other)                          :9675
##
##                               branch    p_id    accuracy    react_time
## edible_fruit.n.01                :1650    1      : 396    Min.    :0.0000    Min.    : 29
## musical_instrument.n.01:2100      2      : 396    1st Qu.:1.0000    1st Qu.: 457
## clothing.n.01                     :2250    3      : 396    Median :1.0000    Median : 563
## hand_tool.n.01                    :1950    4      : 396    Mean   :0.8851    Mean   : 648
## furniture.n.01                    :1950    5      : 396    3rd Qu.:1.0000    3rd Qu.: 736
##                                   6      : 396    Max.   :1.0000    Max.   :8148
##                                   (Other):7524
```

The counts of all stimuli, levels and participant IDs should be the same within their group. If some are higher than others, you might have duplicates and need to take another look at the input data before the evaluation. Here, you can also already see the overall mean accuracy of the participant answers, as well as the mean and median reaction times.

1.2 Reaction time and Mistake Overview

The following box plots in Fig.1 give a short overview of the reaction time means and quantiles with regards to the category levels and also to the stimulus type. A closer look into reaction times can be found in section 4.

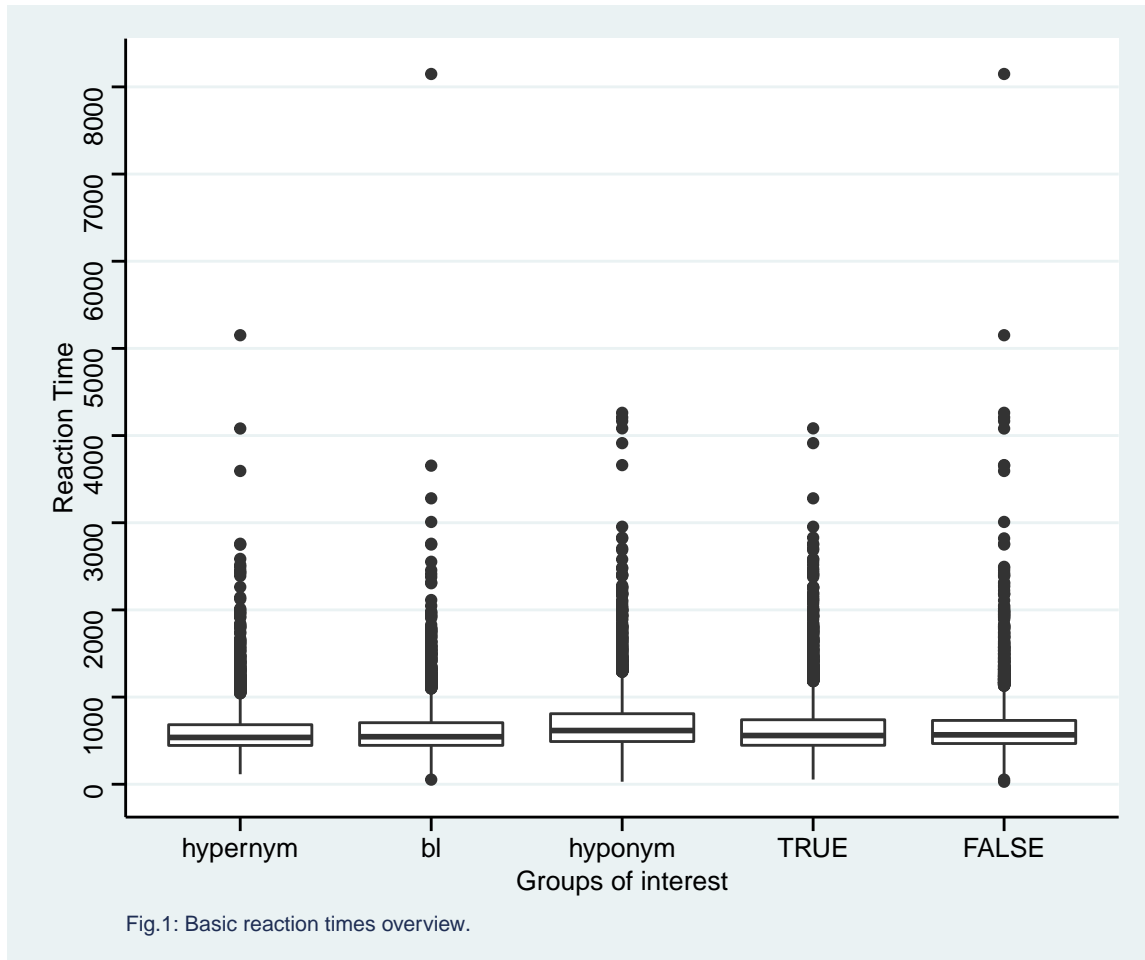


Figure 2 is meant to give an overview of the amount of mistakes made by participants with regard to the main category branches. If the branches are balanced (same amount of categories in each group), then the columns in the plot should have the same height.



Fig.2. Proportions of correct to incorrect participant answers.

2 Mistake Analysis w.r.t. Participant

This section is concerned with mistakes in the experiment from the participant point of view. It helps to determine the performance of participants. It can be used to find participants that should not be rewarded for not doing the job seriously (if it was a condition in the job description) or to reward bonuses (if there are any) to high-performing participants. The per-participant mean reaction times and accuracies are shown.

```
##      p_id      mean_rt      mean_acc
## 1      : 1    Min.    : 414.1    Min.    :0.7929
## 2      : 1    1st Qu.: 534.2    1st Qu.:0.8662
## 3      : 1    Median : 644.3    Median :0.8914
## 4      : 1    Mean    : 648.0    Mean    :0.8851
## 5      : 1    3rd Qu.: 715.5    3rd Qu.:0.9167
## 6      : 1    Max.    :1083.5    Max.    :0.9520
## (Other):19
```

Figure 3 compares the mean accuracies (% on y axis) and reaction times (values in columns) of all participants. It only shows the number of the participant, their actual ID from the job must be looked up manually. Research shows that an accuracy of up to 85% is to be expected [1], anything below could be considered unserious. Look at the following evaluations to verify such a suspicion.

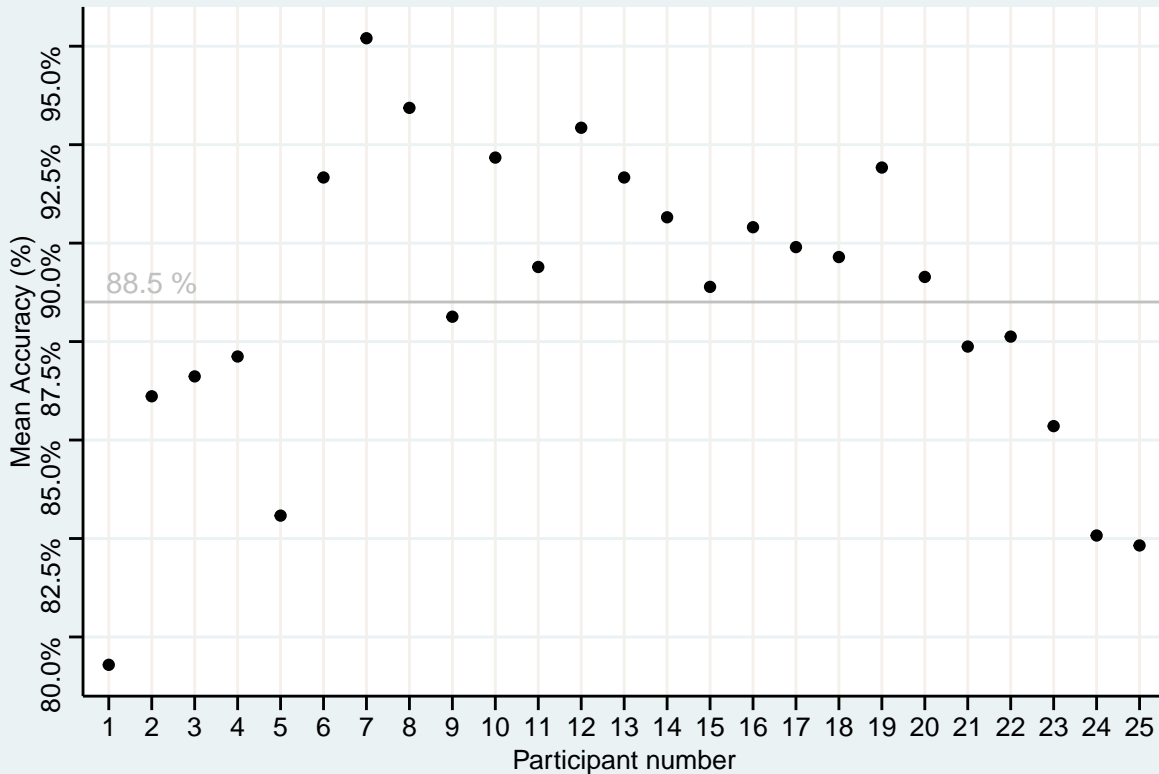
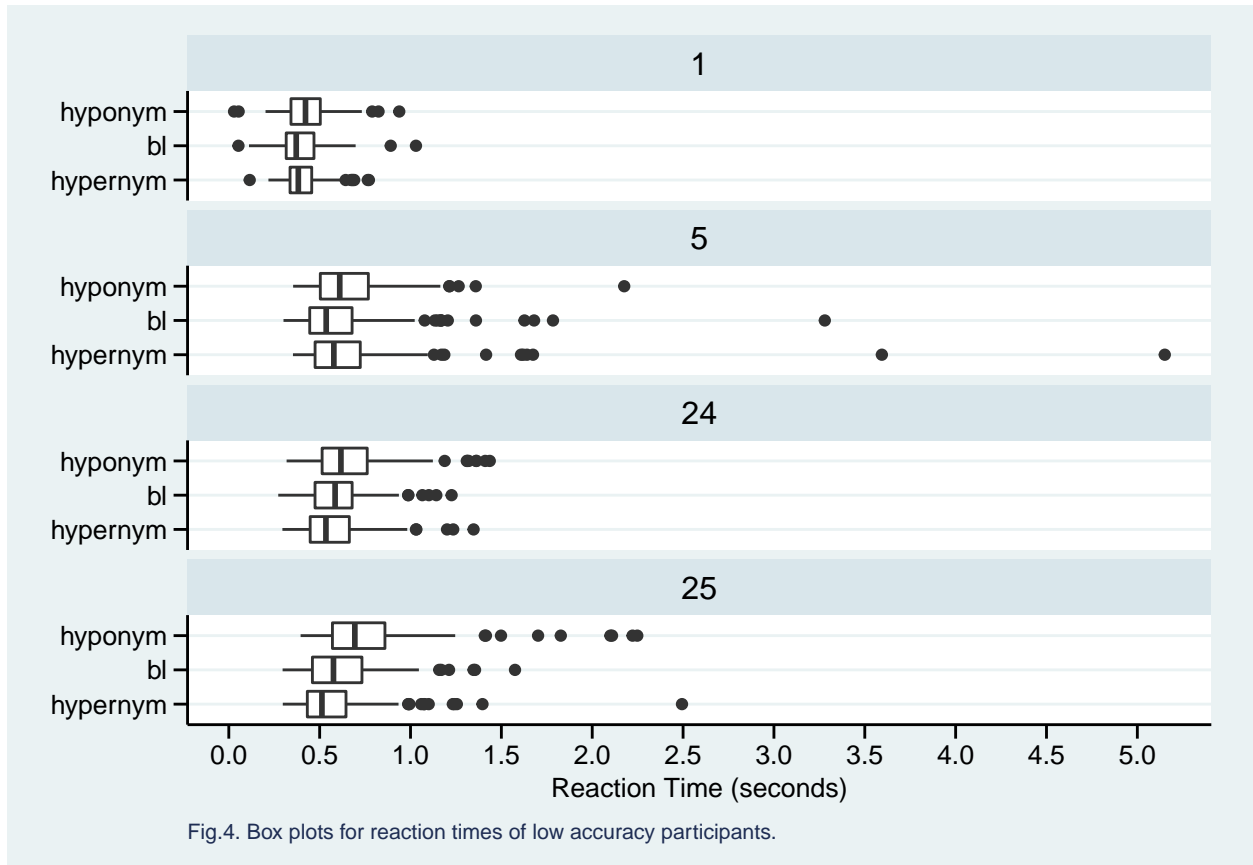


Fig.3. Mean accuracies of all participants.

The following plot (Fig.4) shows box plots for participant's reaction times on stimuli. It is restricted to those participants that show a sub-standard accuracy (<85%). If reaction times are all extremely short, then the participant did not do the task appropriately. Generally, the subordinate reaction times should be slower than the rest. If the range of reaction times is great, then the participant might have struggled with the category names. Outliers above the 3rd quantile indicate that the participant got stuck on those

stimuli / thought longer about them. Outliers below the 1st quantile are answers given too fast to have even registered the stimulus.



The final plot (Fig.5) of this section highlights how participants individually experienced the experiment over time. Each path represents one participant's mistake count as time moves on. Time that passed equally fast for each participant (e.g. showing the stimulus, moving on to the next) is not considered. Only time that passed while the participant was able to answer is considered. Thus, this plot does not show moments where participants took a break during the experiment's break-time screens.

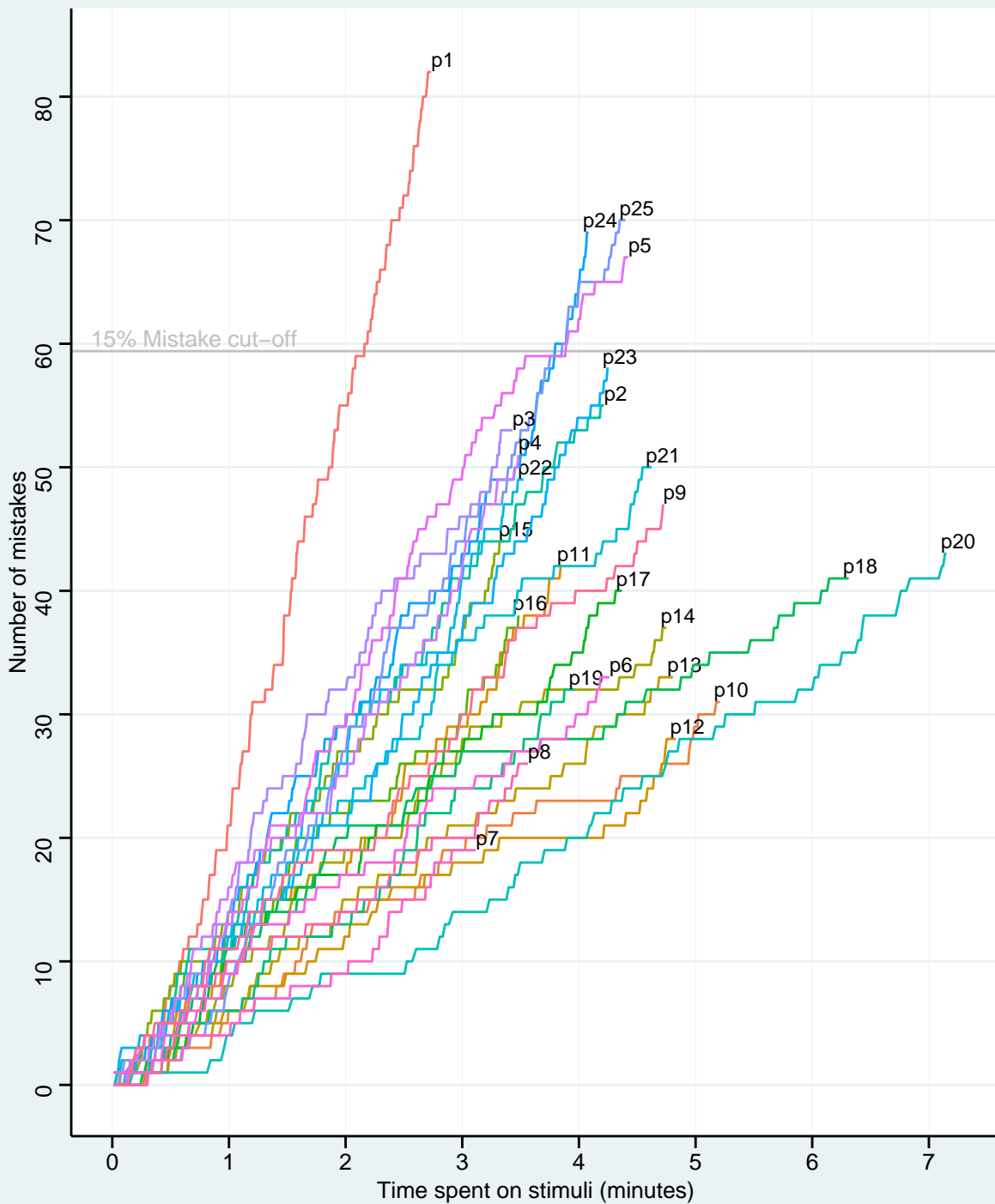
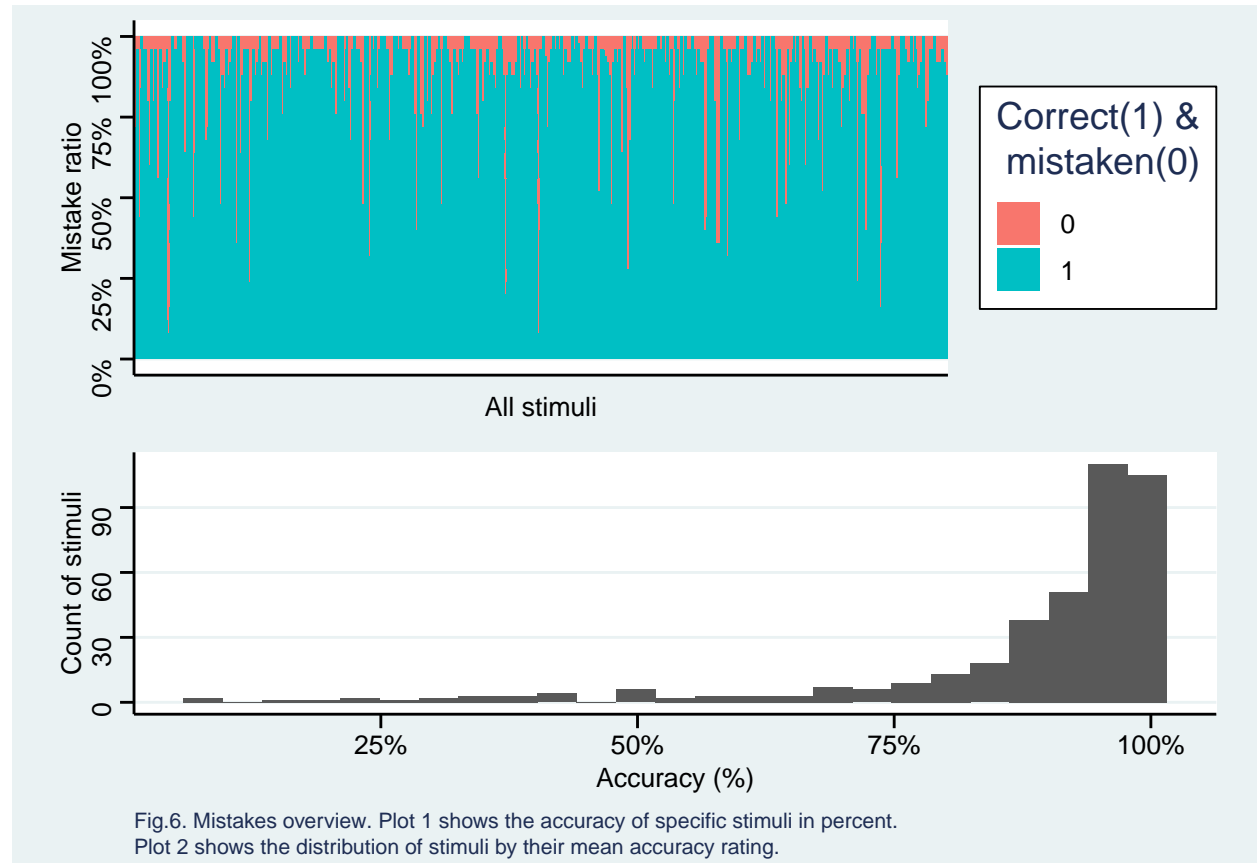


Fig.5. Line graphs for each participant's mistake count over time.

These information should be enough to get a good idea why some participants did better or worse than others. It could also already indicate a possible issue with the experiment procedure or content. The next section will take a closer look at mistakes with respect to the stimuli and their inherent qualities.

3 Mistake Analysis w.r.t. Stimulus

In the following overview (Fig.6), you can see the percentages of mistakes to correct answers of all stimuli. If all participants have seen the stimuli in the same order, then the columns will be in order of how stimuli have been originally ordered.



The following sections try to shine light on what caused mistakes with rights to one specific component or aspect of the stimulus.

3.1 Stimuli

Figure 7 shows the stimuli (image + label) that received the most erroneous answers. The cutoff is made at 50%, when half the participants had a mistake.

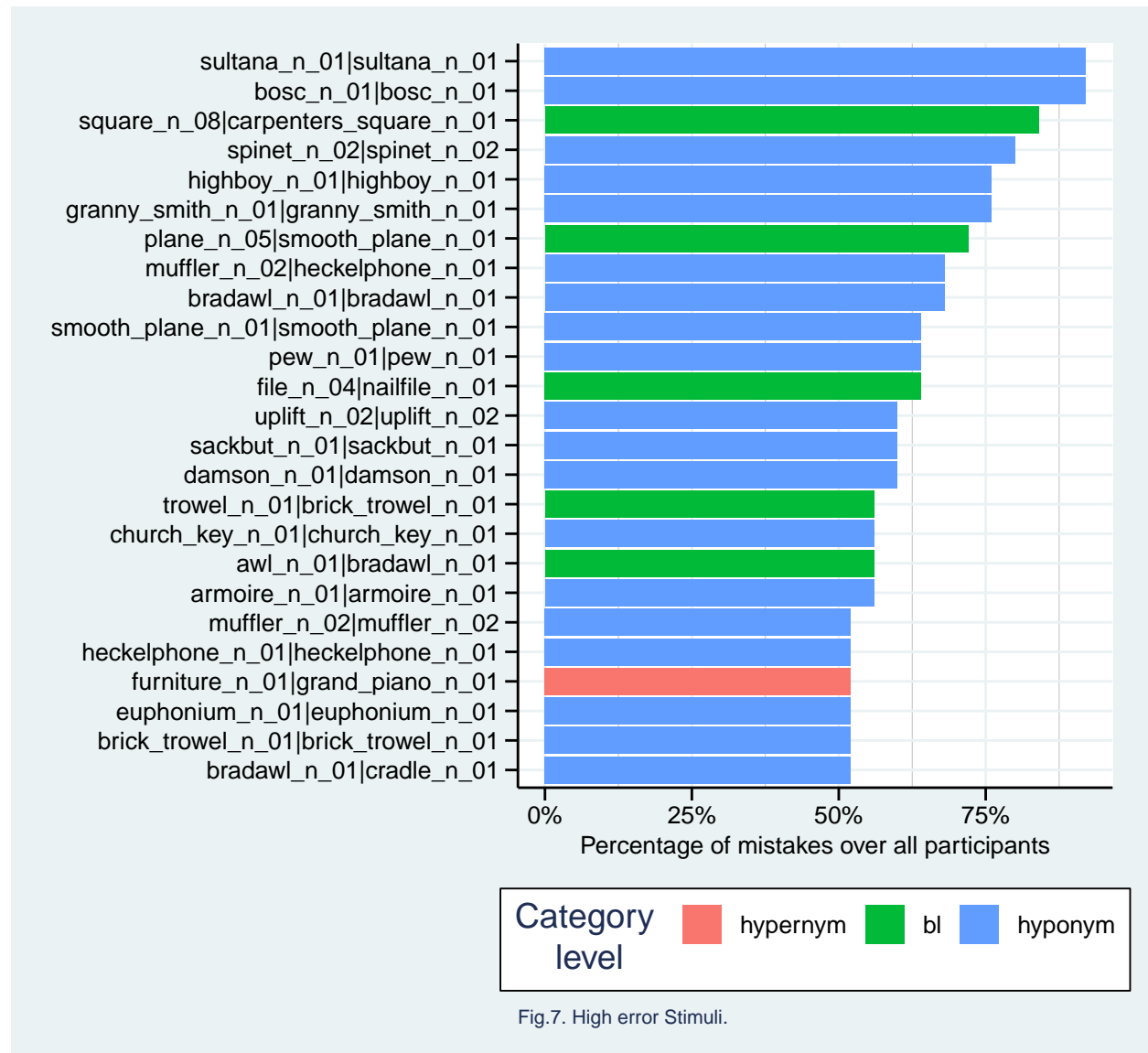
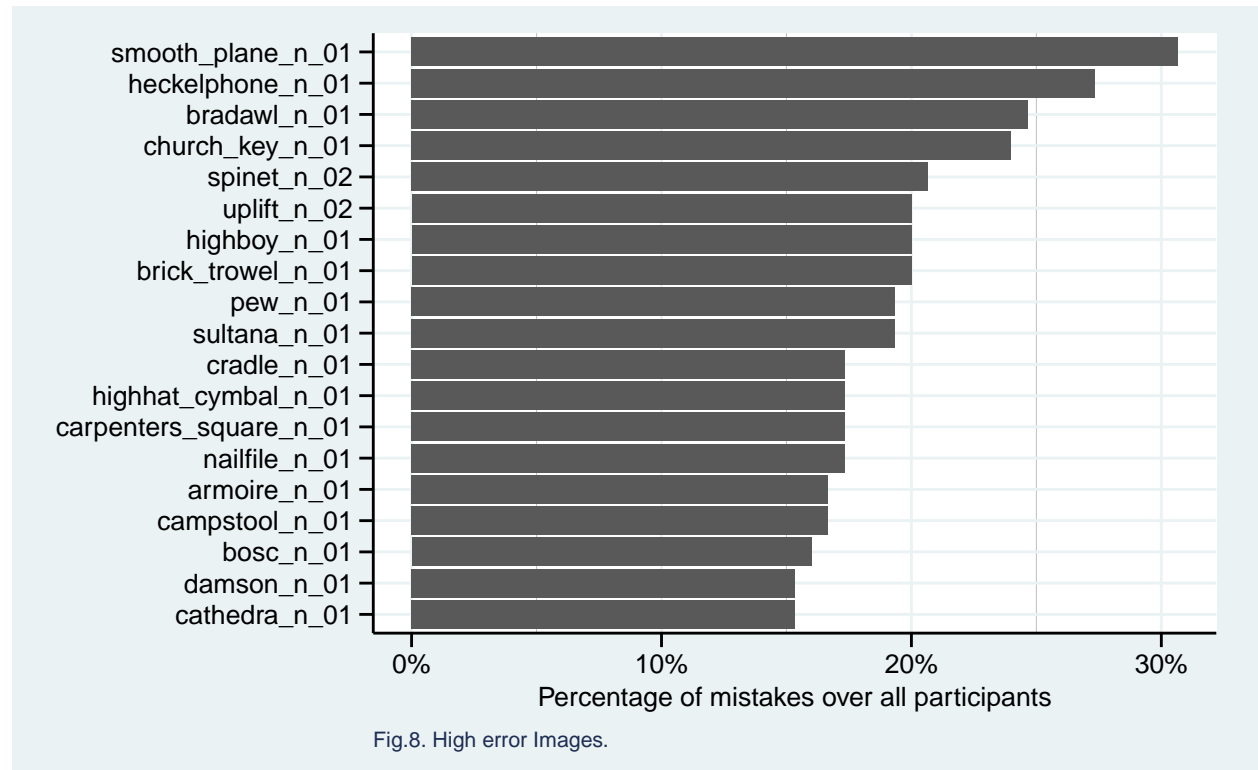


Fig.7. High error Stimuli.

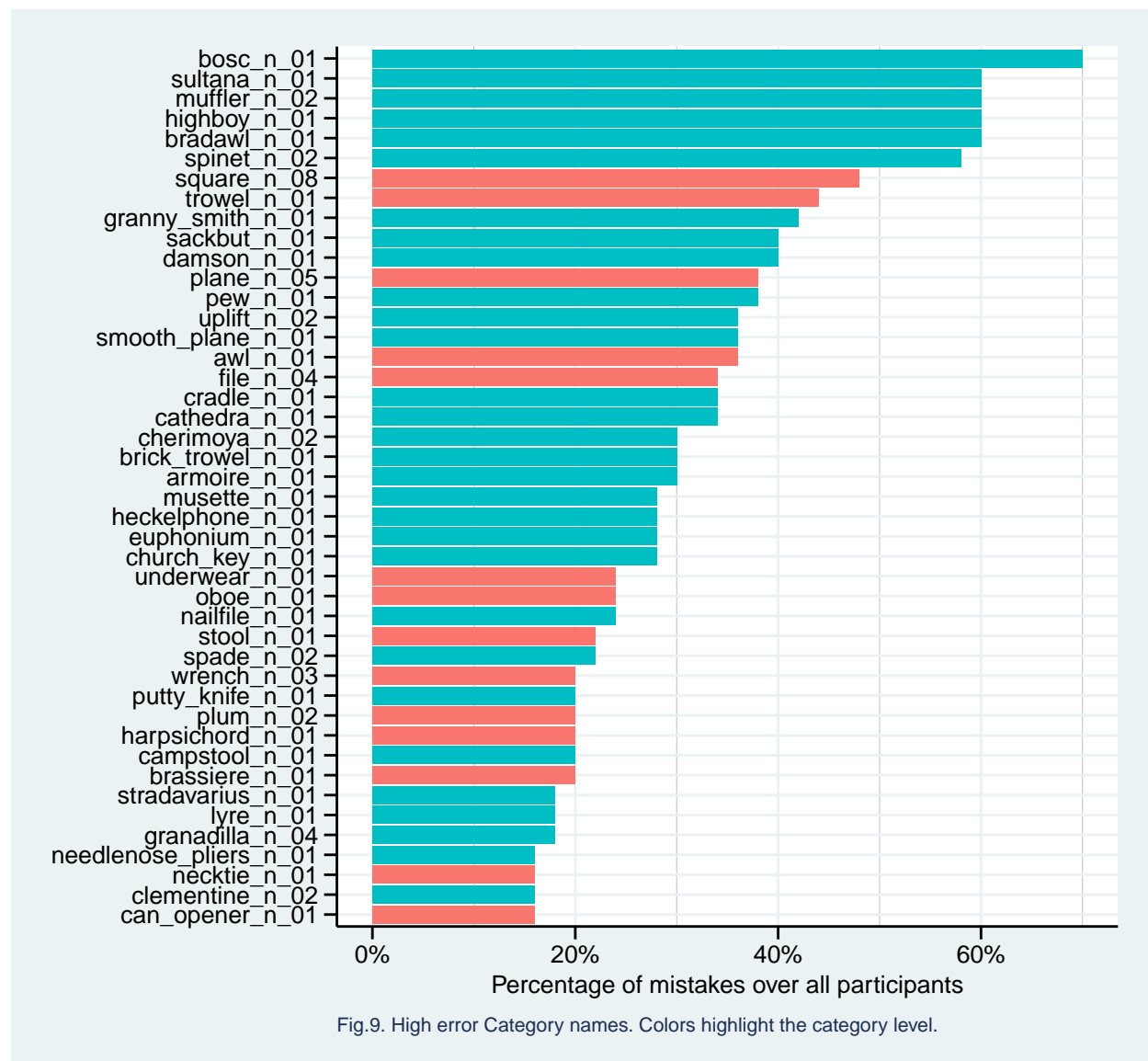
3.2 Image

Figure 8 shows the images that received the most erroneous answers. The cutoff is made at 15%.



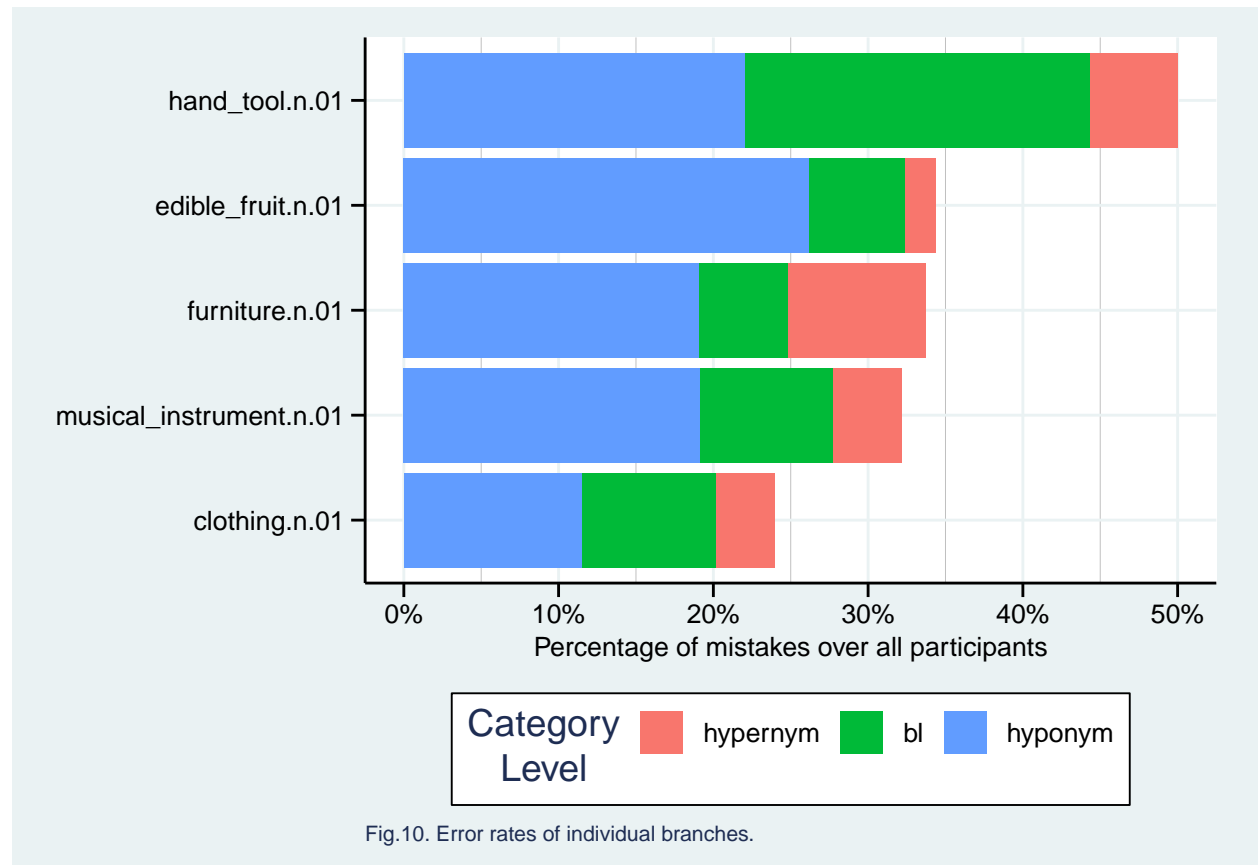
3.3 Category Name

Figure 9 shows the category names that received the most erroneous answers. The cutoff is made at 15%.



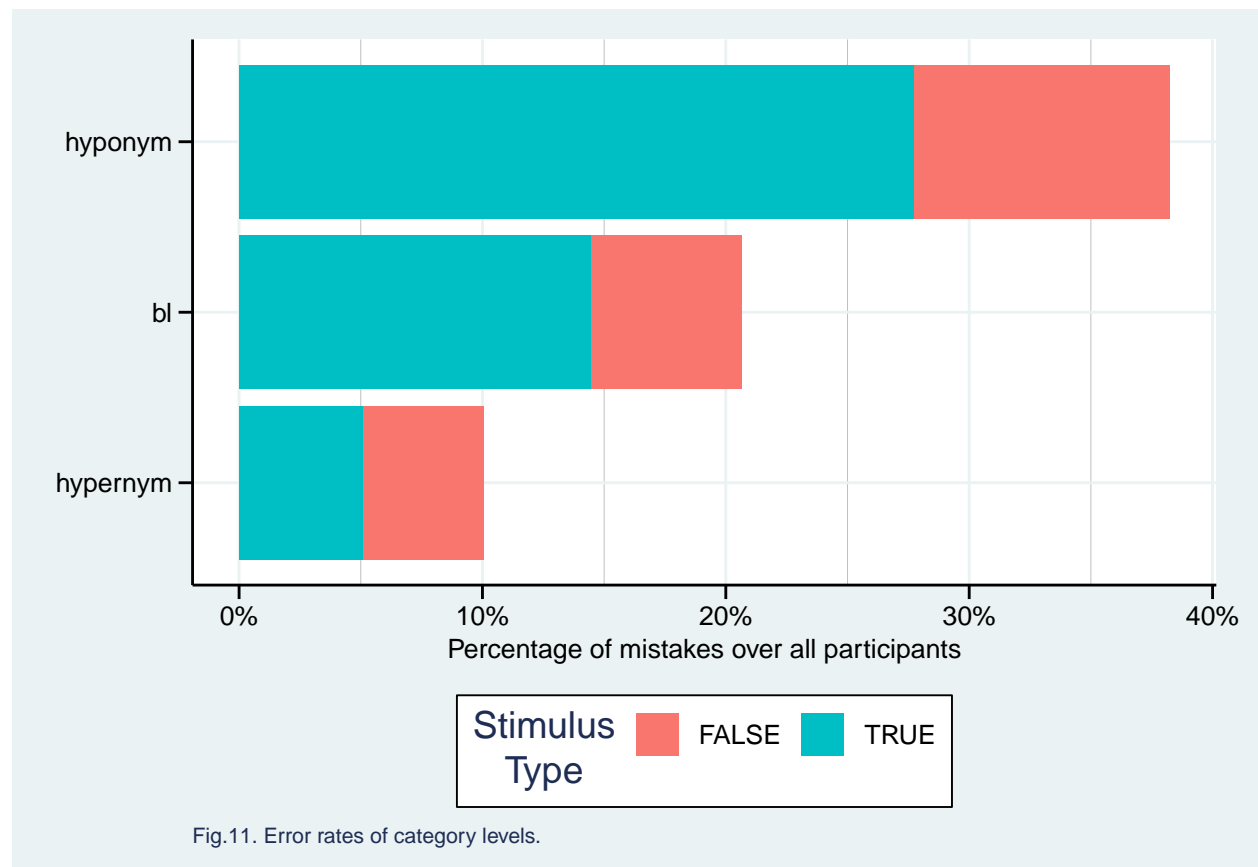
3.4 Branch

Figure 10 shows the mistake percentage per branch. The columns are colored to highlight the share of each category level within the respective branches.



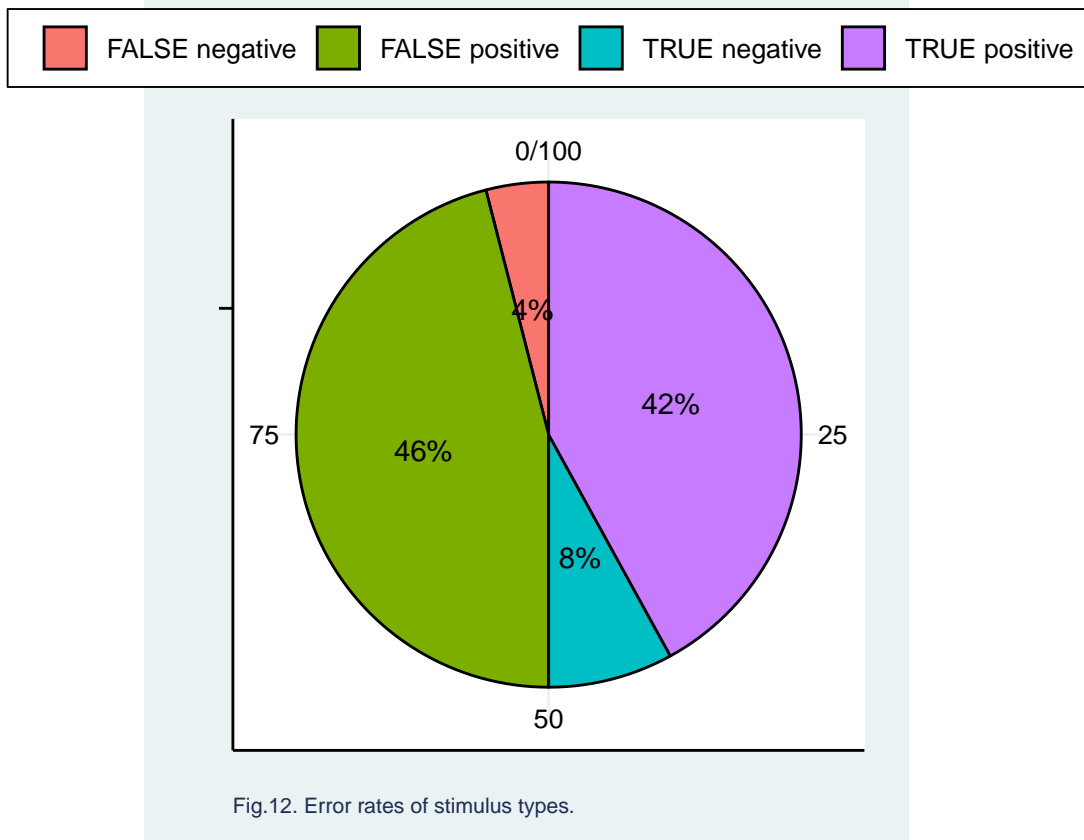
3.5 Level

Figure 11 shows the mistake percentage per category level. The columns are colored to highlight the share of both stimulus types.



3.6 Stimulus Type

Figure 12 shows a pie chart of True/False stimulus type's correct and incorrect answer rates.



4 Component significance analysis in terms of error rates

This section calculates the one-way ANOVA and two-way ANOVA of each factor (stimulus components) to determine significant differences between groups. Then, they are pitted against each other in an AIC table to determine the best-fitting statistical model among the ANOVA models.

4.1 One-Way

```
aov_branch <- aov(err ~ branch, data=aov_data)
aov_lvl <- aov(err ~ lvl, data=aov_data)
aov_c_name <- aov(err ~ c_name, data=aov_data)
aov_img <- aov(err ~ img, data=aov_data)
aov_stim <- aov(err ~ stimulus, data=aov_data)
aov_stype <- aov(err ~ stim_type, data=aov_data)
```

4.1.1 Branch

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## branch         4  0.326  0.08162   2.904 0.0217 *
## Residuals     388 10.907  0.02811
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.1.2 Level

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lvl           2  1.323  0.6616   26.04 2.44e-11 ***
## Residuals    390  9.910  0.0254
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.1.3 Category name

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## c_name       136  7.068  0.05197   3.194 6.17e-16 ***
## Residuals    256  4.165  0.01627
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.1.4 Image

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## img          65  1.589  0.02445   0.829  0.819
## Residuals    327  9.644  0.02949
```

4.1.5 Stimulus

```
##              Df Sum Sq Mean Sq
## stimulus     392 11.23  0.02866
```

4.1.6 Stimulus type

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## stim_type      1  0.701  0.7007    26.01 5.3e-07 ***
## Residuals    391 10.532  0.0269
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.1.7 AIC

```
## Warning: package 'AICcmodavg' was built under R version 4.2.1
```

```
##
## Model selection based on AICc:
##
##           K      AICc Delta_AICc AICcWt Cum.Wt      LL
## Branch      6 -281.18      Inf     NaN     NaN 146.70
## Level       4 -322.95      Inf     NaN     NaN 165.53
## Cat. Name  138 -244.67      Inf     NaN     NaN 335.85
## Image      67 -179.72      Inf     NaN     NaN 170.88
## Stim. Type   3 -301.06      Inf     NaN     NaN 153.56
## Stimulus   394    -Inf      NaN     NaN     NaN    Inf
```


4.2 Two-way

In the first iteration of this evaluation, the models with the most significant differences within their groups were those for stimulus type, category level, category name and branch (lower significance). The AIC table determined, in this order, that the ANOVA models for category level, stimulus type and branch showed the best fit. The following subsections compare ANOVA of pairs and triples of these components in another AIC table.

```
aov_lvl_type <- aov(err ~ lvl + stim_type, data=aov_data)
aov_lvl_branch <- aov(err ~ lvl + branch, data=aov_data)
aov_type_branch <- aov(err ~ stim_type + branch, data=aov_data)
aov_lvl_type_branch <- aov(err ~ lvl + stim_type + branch, data=aov_data)
```

4.2.1 Level + Type

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lvl          2  1.323   0.6616   28.00 4.38e-12 ***
## stim_type    1  0.717   0.7173   30.36 6.56e-08 ***
## Residuals   389  9.193   0.0236
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.2.2 Level + Branch

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lvl          2  1.323   0.6616  26.642 1.45e-11 ***
## branch       4  0.325   0.0813   3.275  0.0117 *
## Residuals   386  9.585   0.0248
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.2.3 Type + Branch

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## stim_type    1  0.701   0.7007  26.567 4.07e-07 ***
## branch       4  0.326   0.0814   3.087  0.016 *
## Residuals   387 10.207   0.0264
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.2.4 Level + Type + Branch

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## lvl          2  1.323   0.6616  28.721 2.37e-12 ***
## stim_type    1  0.717   0.7173  31.142 4.53e-08 ***
## branch       4  0.324   0.0811   3.522  0.00772 **
## Residuals   385  8.868   0.0230
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

4.2.5 AIC

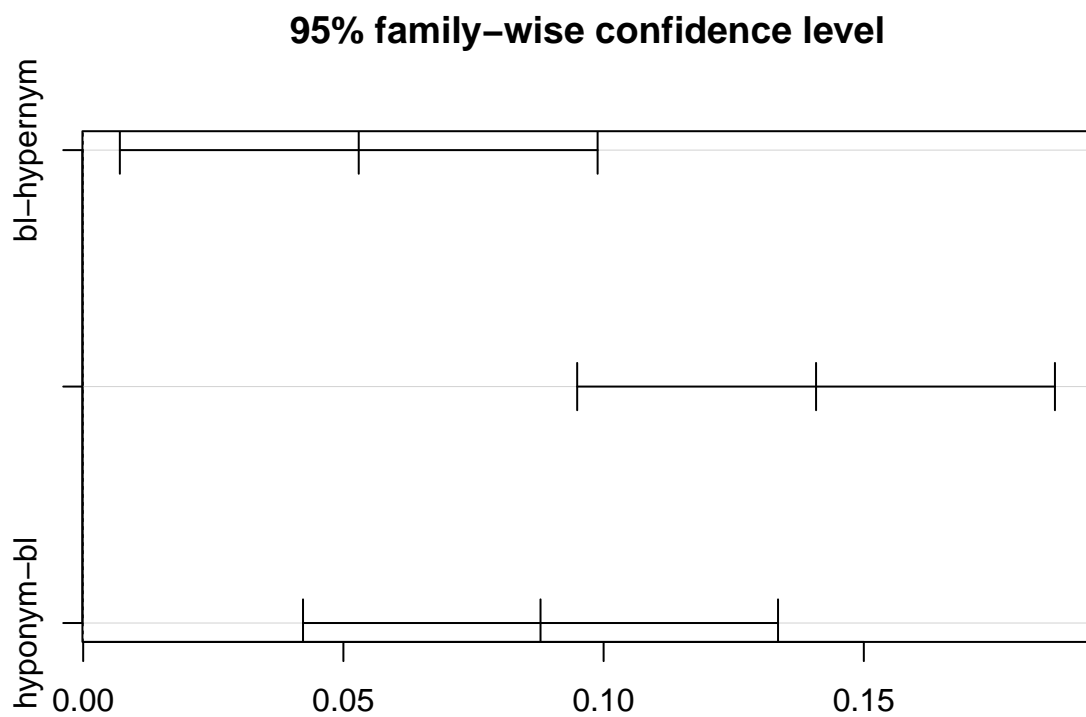
```
##
## Model selection based on AICc:
##
##           K      AICc Delta_AICc AICcWt Cum.Wt      LL
## Level & Type & Branch 9 -356.24      0.00  0.95  0.95 187.35
## Level & Type          5 -350.43      5.81  0.05  1.00 180.29
## Level & Branch        8 -327.80     28.44  0.00  1.00 172.09
## Type & Branch         7 -305.17     51.07  0.00  1.00 159.73
```

A Tukey test can be used to measure the differences between group-member pairings.

```
tukey <- TukeyHSD(aov_lvl_branch)
print(tukey)

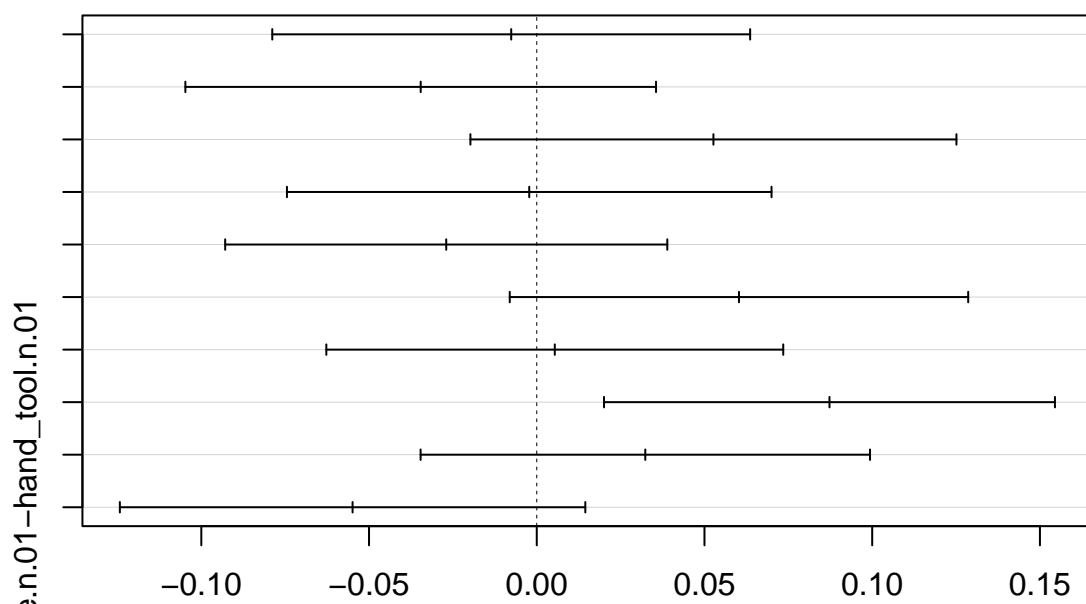
##      Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = err ~ lvl + branch, data = aov_data)
##
## $lvl
##           diff           lwr           upr      p adj
## bl-hypernym      0.05294574 0.007045087 0.09884639 0.0189998
## hyponym-hypernym 0.14082452 0.094923875 0.18672517 0.0000000
## hyponym-bl       0.08787879 0.042242698 0.13351488 0.0000234
##
## $branch
##                               diff           lwr           upr
## musical_instrument.n.01-edible_fruit.n.01 -0.007612805 -0.078839619 0.06361401
## clothing.n.01-edible_fruit.n.01           -0.034597029 -0.104752295 0.03555824
## hand_tool.n.01-edible_fruit.n.01           0.052667661 -0.019777951 0.12511327
## furniture.n.01-edible_fruit.n.01           -0.002237762 -0.074468720 0.06999320
## clothing.n.01-musical_instrument.n.01      -0.026984224 -0.092885158 0.03891671
## hand_tool.n.01-musical_instrument.n.01      0.060280466 -0.008053555 0.12861449
## furniture.n.01-musical_instrument.n.01      0.005375043 -0.062731367 0.07348145
## hand_tool.n.01-clothing.n.01                0.087264690  0.020048316 0.15448106
## furniture.n.01-clothing.n.01                0.032359266 -0.034625698 0.09934423
## furniture.n.01-hand_tool.n.01              -0.054905423 -0.124285469 0.01447462
##
##                               p adj
## musical_instrument.n.01-edible_fruit.n.01 0.9983894
## clothing.n.01-edible_fruit.n.01           0.6590137
## hand_tool.n.01-edible_fruit.n.01           0.2715710
## furniture.n.01-edible_fruit.n.01           0.9999882
## clothing.n.01-musical_instrument.n.01      0.7947523
## hand_tool.n.01-musical_instrument.n.01      0.1127500
## furniture.n.01-musical_instrument.n.01      0.9995124
## hand_tool.n.01-clothing.n.01               0.0038262
## furniture.n.01-clothing.n.01               0.6764229
## furniture.n.01-hand_tool.n.01              0.1937813
```

```
plot(tukey, sub="Fig.13. Tukey Confidence intervals for level and branch.")
```



Differences in mean levels of lvi
Fig.13. Tukey Confidence intervals for level and branch.

95% family-wise confidence level



Differences in mean levels of branch
Fig.13. Tukey Confidence intervals for level and branch.

5 Analysis of reaction times

This section is concerned with the reaction times recorded during the experiment. The per-participant reaction times have already been illustrated in Fig.3.

5.1 Overview

To give an overview, Fig.14 illustrates the distribution of stimuli over the range of reaction times. The columns for the three distinct category levels sit next to each other to facilitate the identification of possible differences. Fig.15 shows box plots for reaction times within the three category levels, one plot per stimulus type. Fig.15 is also a visual representation of table 1, which is the same table Rosch[2] created.

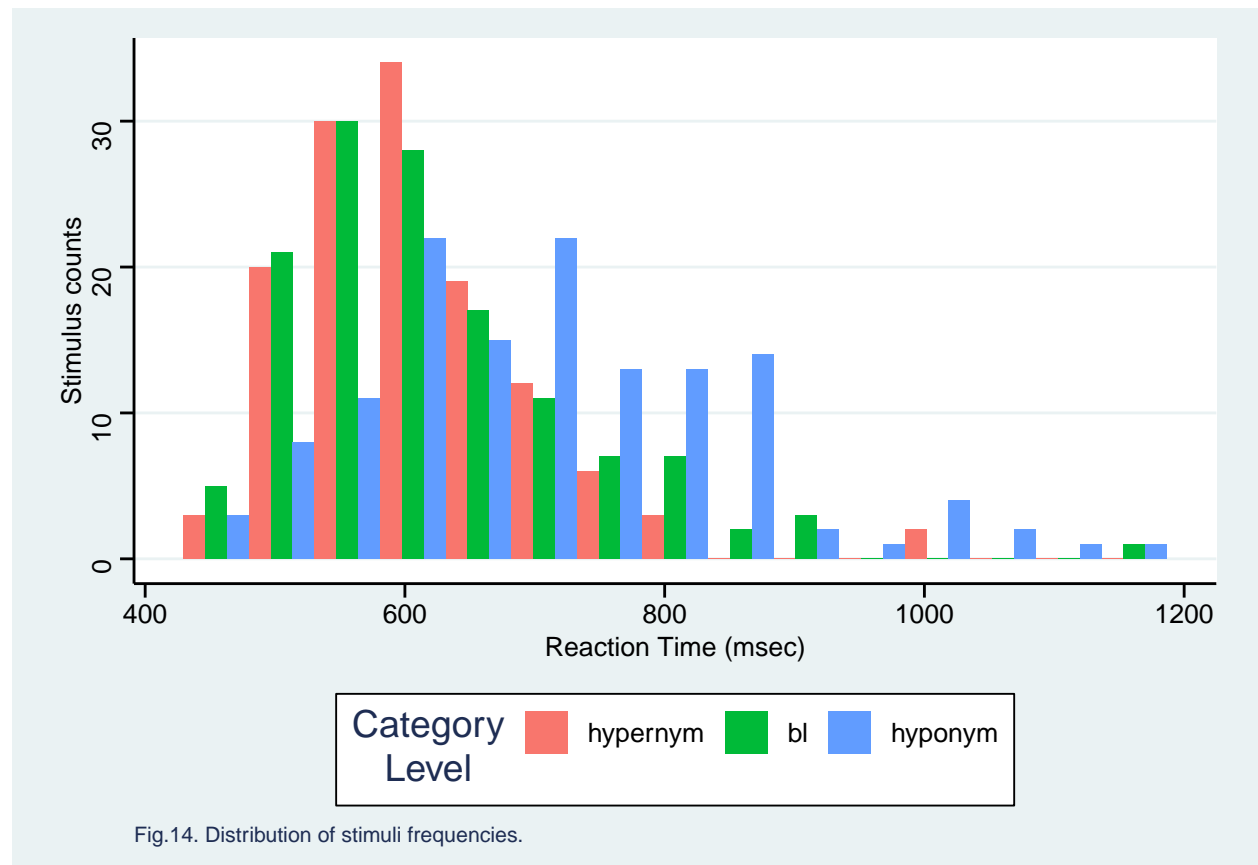


Fig.14. Distribution of stimuli frequencies.

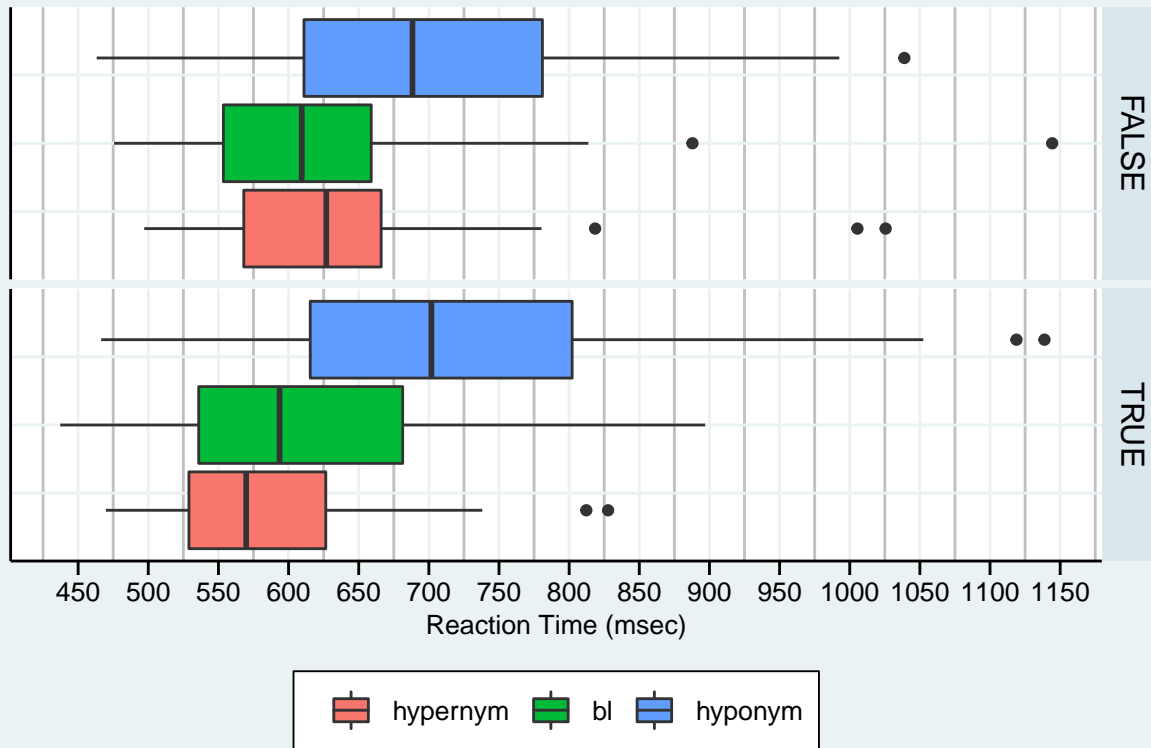


Fig.15 Boxplot representation of reaction time means over all stimuli

##	Superordinate	Basic Level	Subordinate
## F	635.6127	623.7824	696.9606
## T	585.9824	621.8145	723.8018

Table 1.: Matrix showing the mean reaction times at different category levels and stimuli types.

5.2 Analysis

This section applies the two-way ANOVA model on the measurements. The category level (between-subject fixed effect) and the category names (random variable) are used as independent variables.

5.2.1 True type stimuli

Anova summary:

```
summary.aov(roschAOV_true)
```

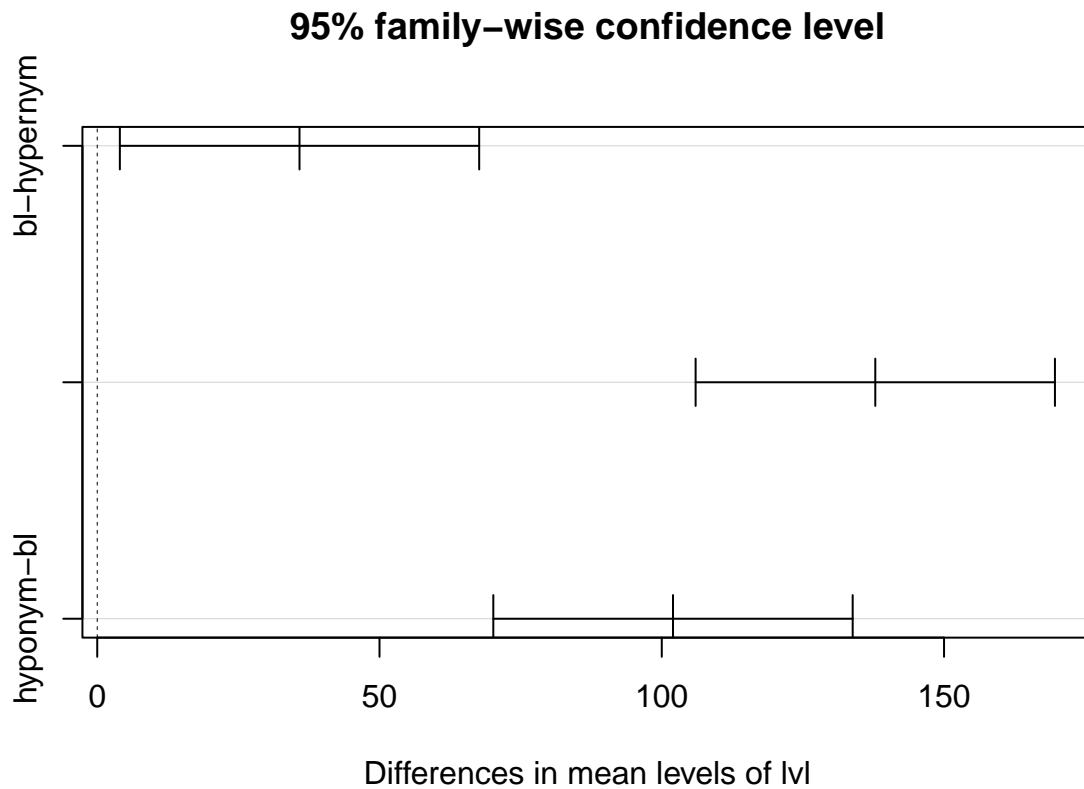
```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## lvl         2  674950   337475   58.261 7.09e-15 ***
## c_name     134 2554301    19062    3.291 3.98e-07 ***
## Residuals   61  353341     5792
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Tukey test:

```
tukey_r_true
```

```
## Tukey multiple comparisons of means
## 95% family-wise confidence level
##
## Fit: aov(formula = rt_mean ~ lvl + c_name, data = r_true)
##
## $lvl
##          diff      lwr      upr    p adj
## bl-hypernym   35.83212   4.00569  67.65855 0.0237145
## hyponym-hypernym 137.81939 105.99296 169.64583 0.0000000
## hyponym-bl    101.98727  70.16084 133.81370 0.0000000
```

```
plot(tukey_r_true)
```



5.2.2 False type stimuli

Anova summary:

```
summary.aov(roschAOV_false)
```

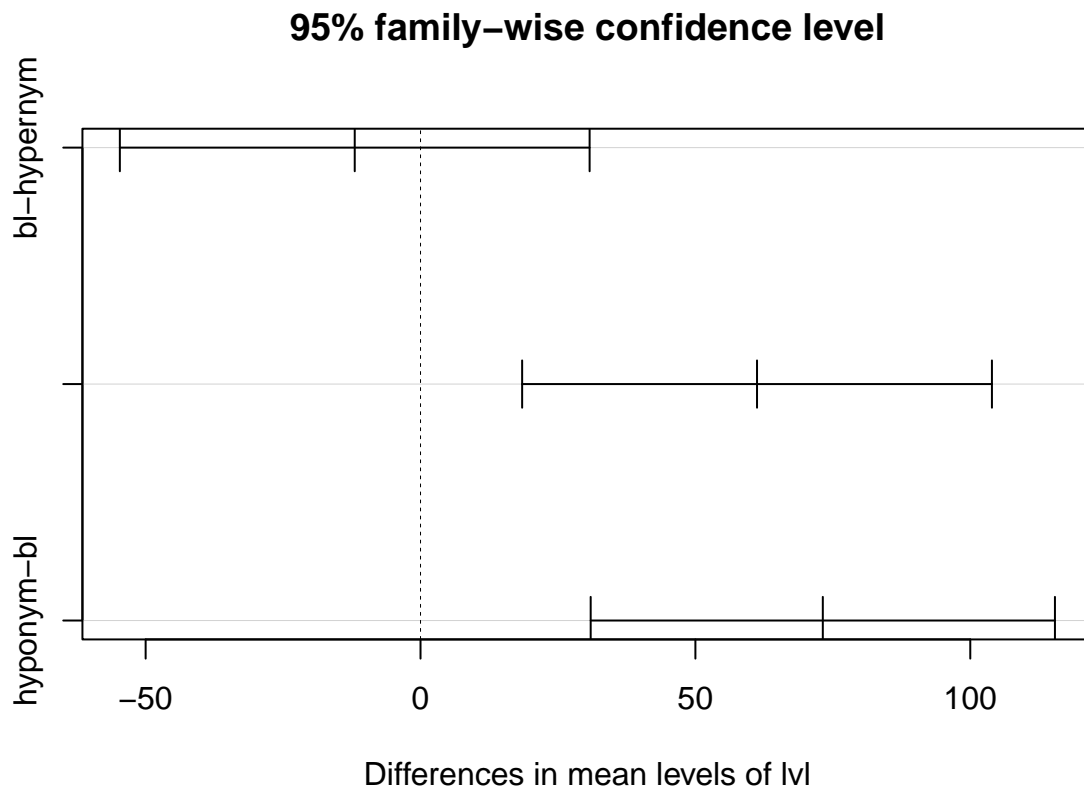
```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## lvl         2  202575   101288    9.959 0.000191 ***
## c_name     134 1777139    13262    1.304 0.127242
## Residuals   58  589892    10171
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Tukey test:

```
tukey_r_false
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = rt_mean ~ lvl + c_name, data = r_false)
##
## $lvl
##           diff          lwr          upr      p adj
## bl-hypernym   -11.96488 -54.69129  30.76153 0.7797189
## hyponym-hypernym 61.21330  18.48689 103.93972 0.0030132
## hyponym-bl      73.17818  30.95151 115.40485 0.0003015
```

```
plot(tukey_r_false)
```



6 B.L. Determination

Ultimately, the experiment is conducted to determine the basic level category name from a selection of three candidates. These triples form a path through one of the WordNet branches. We only consider a single element in this path basic level. Basic levelness can only be inferred from reaction times of the true type stimuli. The false type stimuli reactions say more about the image than about the category names.

6.1 Simple rating

In this section, we will see a simple implementation to declare which category name from the triple is the basic level (super- and subordinates can be inferred). For each triple, the category name with the fastest mean reaction time is chosen as basic level. Below, you can see the top nine results of the evaluation. Thereafter follows a confusion matrix of expected and predicted levels.

```
## # A tibble: 9 x 3
## # Groups:   lvl, c_name [7]
##   c_name          lvl proj_lvl
##   <chr>          <fct> <fct>
## 1 clothing_n_01    super bl
## 2 sweater_n_01     bl    sub
## 3 turtleneck_n_01 sub    sub
## 4 clothing_n_01    super bl
## 5 shirt_n_01       bl    sub
## 6 tank_top_n_01    sub    sub
## 7 clothing_n_01    super super
## 8 necktie_n_01     bl    bl
## 9 bow_tie_n_01     sub    sub

## Warning: package 'caret' was built under R version 4.2.1

## Loading required package: lattice

##
## Attaching package: 'caret'

## The following object is masked from 'package:purrr':
##
##   lift

## Confusion Matrix and Statistics
##
##           Reference
## Prediction bl sub super
##      bl    28  4   34
##      sub   38 62    0
##      super  0  0   32
##
## Overall Statistics
##
##           Accuracy : 0.6162
##           95% CI : (0.5446, 0.6842)
```

```

##      No Information Rate : 0.3333
##      P-Value [Acc > NIR] : 3.927e-16
##
##              Kappa : 0.4242
##
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: bl Class: sub Class: super
## Sensitivity          0.4242      0.9394      0.4848
## Specificity          0.7121      0.7121      1.0000
## Pos Pred Value       0.4242      0.6200      1.0000
## Neg Pred Value       0.7121      0.9592      0.7952
## Prevalence           0.3333      0.3333      0.3333
## Detection Rate       0.1414      0.3131      0.1616
## Detection Prevalence 0.3333      0.5051      0.1616
## Balanced Accuracy     0.5682      0.8258      0.7424

```

6.2 Probabilistic rating

The basic level is not an inherent quality of a word itself, it is a quality experienced by humans using the words. Everyone experiences language slightly differently, wherefore basic-levelness should rather be expressed as a probability. Then it can be said, this word has a higher probability of being at the basic level than its hypernym or hyponym.

6.3 Probabilistic rating with error

Stimuli that gathered many True negative responses, but few False positives, might give insight into basic levelness. The participant knew the object, which is why they could correctly determine a label as wrong, but they made many mistakes when confronted with the lower level category name. It could mean that they simply did not know the word, which would considerably lower the probability of it being basic level.

7 References

- [1] Gagné, N., & Franzen, L., 2021, <https://doi.org/10.31234/osf.io/nt67j>
- [2] E. Rosch, C. B. Mervis, W. D. Gray, D. M. Johnson, and P. Boyes-Braem, “Basic objects in natural categories,” *Cognitive Psychology*, vol. 8, no. 3, pp. 382–439, 1976.