# rnnpsycholing Japanese NPI (no embedded sentences)

*Takashi Morita*

## 目次

## Introduction

We are looking for a 2x2 interaction of:

- presence vs. absence of the Japanese NPI *shika* (しか)
- affirmativeness vs. negativeness of main verb

for each of the three grammatical cases (TOP, ACC, DAT) of the *shika*-attached NP.

e.g.

- TOP
  - 渡辺-{しか, は} 家族-に 手紙-を {渡した, 渡さなかった}。
  - Watanabe-{*shika*,TOP} family-DAT letter-ACC {came, didn't come}.
  - 'Only Watanabe handed letters to his family.'
  - 'Watanabe handed/didn't hand letters to his family.'
- ACC
  - 渡辺-は 家族-に 手紙-{しか, を} {渡した, 渡さなかった}。
  - Watanabe-TOP family-DAT letter-{*shika*,ACC} {came, didn't.come.}
  - 'Watanabe handed only letters to his family.'
  - 'Watanabe handed/didn't hand letters to his family.'
- DAT
  - 渡辺-は 家族-{に-しか, に} 手紙-を {渡した, 渡さなかった}。
  - Watanabe-TOP family-{DAT-*shika*,DAT} letter-ACC {came, didn't.come.}
  - 'Watanabe handed letters only to his family.'

– 'Watanabe handed/didn't hand letters to his family.'

Why is this interesting?

- A grammatical sentence with *shika* must have a negative verb.
- Affirmative verbs would show significant increase in surprisal when *shika* precedes compared with its absence.

## Methods

For each pair $i$ of sentences with vs. without *shika*, we look at their difference in surprisal of the verb ($V$) region.

$$D_i := S(V_i \mid \texttt{shika}) - S(V_i \mid \texttt{no-shika})$$

$$S(r) := -\log_2 P(r)$$

And we perform a statistical analysis and check if the affirmativeness vs. negativeness of the verb have an effect on the surprisal difference $D$.

## Load data

```
rm(list = ls())
library(tidyverse)
library(brms)
library(lme4)
library(lmerTest)
library(plotrix)



REGIONS = c('prefix', 'V', 'end')



token_based_data_path = 'jp_shika_test_sentences_unembedded_surprisal-per-token.tsv'
data_token_based = read_tsv(token_based_data_path)

## Parsed with column specification:
## cols(
##   sent_index = col_integer(),
##   token_index = col_integer(),
##   token = col_character(),
##   region = col_character(),
```

```
##    log_prob = col_double(),
##    shika_case = col_character(),
##    shika_embedded = col_character(),
##    other_v_type = col_character(),
##    shika = col_character(),
##    verb_type = col_character(),
##    surprisal = col_double(),
##    LSTM = col_character()
## )
```

```r
# Fill the initial surprisal by 0.
data_token_based[is.na(data_token_based$surprisal),]$surprisal = 0
data_token_based$region = factor(data_token_based$region, levels=REGIONS)



data_region_based = data_token_based %>%
    group_by(sent_index, region, shika, verb_type, shika_case) %>%
        summarise(surprisal=sum(surprisal)) %>%
        ungroup() %>%
    mutate(
        shika=factor(shika, levels=c("shika", "no-shika")),
        verb_type=factor(verb_type, levels=c("affirmative", "negative")),
        shika_case=factor(shika_case, levels=c("TOP", "ACC", "DAT"))
        )

# Sum coding of the variables.
contrasts(data_region_based$shika) = "contr.sum"
contrasts(data_region_based$verb_type) = "contr.sum"

# Make sure that the dataframe is sorted appropriately.
# First by case
data_region_based = data_region_based[order(data_region_based$shika_case),]
# Second by verb_type (affirmative vs. negative)
data_region_based = data_region_based[order(data_region_based$verb_type),]
# Then by sent_index
data_region_based = data_region_based[order(data_region_based$sent_index),]
```
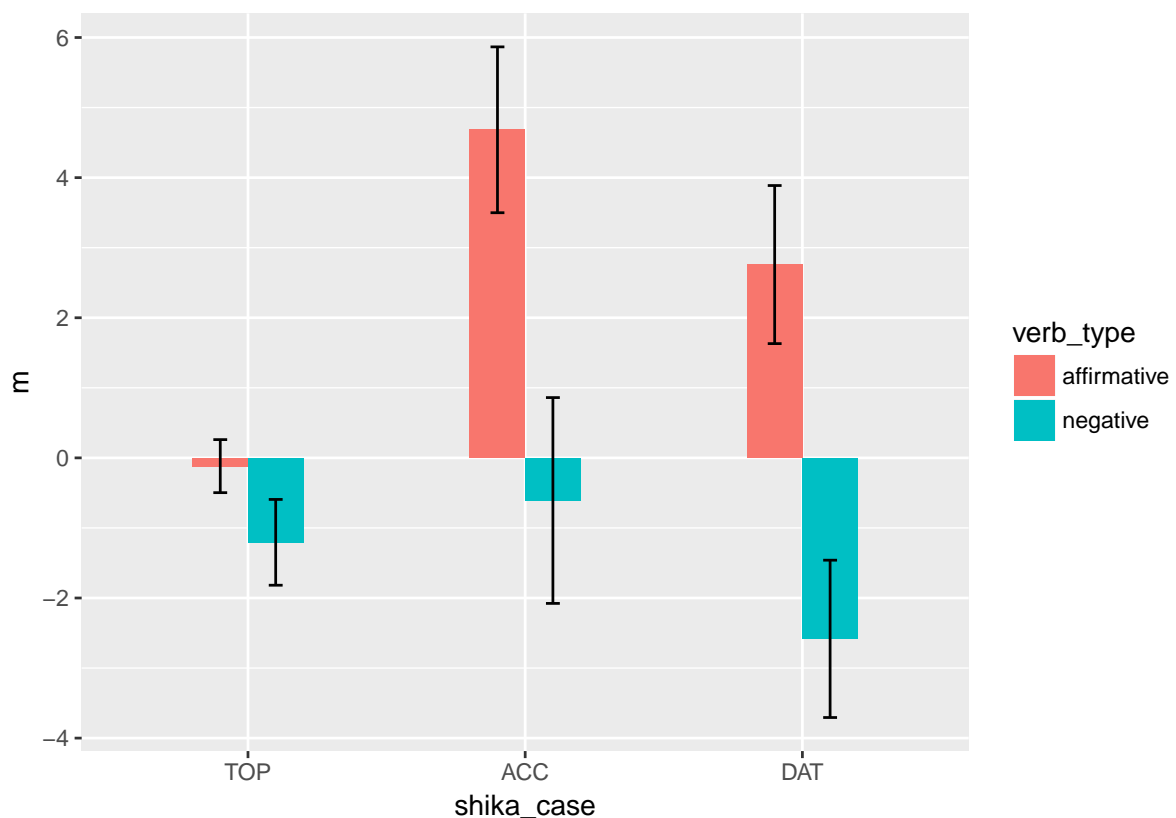
## Visualization

```r
# Focus on the V (verb) region.
data_V = subset(data_region_based, region == 'V')
```

```
# Get difference in surprisal between shika vs. no-shika.
data_V_shika = subset(data_V, shika == 'shika')
data_V_no_shika = subset(data_V, shika == 'no-shika')
data_V_shika$surprisal_diff = data_V_shika$surprisal - data_V_no_shika$surprisal


# Visualize the difference in surprisal increase/dicrease between affirmative vs. negative verbs.
data_V_shika %>%
    group_by(verb_type, shika_case) %>%
    summarise(m=mean(surprisal_diff),
            s=std.error(surprisal_diff),
            upper=m + 1.96*s,
            lower=m - 1.96*s) %>%
    ungroup() %>%
    ggplot(aes(x=shika_case, y=m, ymin=lower, ymax=upper, width=0.4, fill=verb_type)) +
        geom_bar(stat = 'identity', position = "dodge") +
        geom_errorbar(position=position_dodge(0.4), width=.1)
```



- TOP - No visible increase in surprisal of the affirmative verbs. - Visible decrease in surprisal of the negative verbs. - ACC - Greatest increase in surprisal of the affirmative verbs. - Small decrease in surprisal of the negative verbs. - DAT - Visible increase in surprisal of the affirmative verbs. - Greatest decrease in surprisal of the negative verbs.

# Regressions

## TOP

```
sub_data = subset(data_V_shika, shika_case == 'TOP')


m = lmer(
        surprisal_diff
            ~ verb_type
                + (1 | sent_index)
        ,
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: surprisal_diff ~ verb_type + (1 | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 376.7
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.39656 -0.38524 -0.04051  0.46275  2.27091
##
## Random effects:
##  Groups      Name        Variance Std.Dev.
##  sent_index (Intercept) 1.393    1.180
##  Residual               1.843    1.358
## Number of obs: 96, groups:  sent_index, 48
##
## Fixed effects:
##             Estimate Std. Error       df t value Pr(>|t|)
## (Intercept)  -0.6617     0.2196 47.0000  -3.013 0.004159 **
## verb_type1    0.5436     0.1386 47.0000   3.923 0.000283 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##            (Intr)
## verb_type1 0.000
```

- Statistically significant effect of the affirmativeness vs. negativeness of verbs.

## ACC

```
sub_data = subset(data_V_shika, shika_case == 'ACC')


m = lmer(
        surprisal_diff
            ~ verb_type
                + (1 | sent_index)
        ,
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: surprisal_diff ~ verb_type + (1 | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 171.4
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -2.03979 -0.47071  0.05871  0.53181  1.37674
##
## Random effects:
##  Groups      Name        Variance Std.Dev.
##  sent_index (Intercept) 6.79     2.606
##  Residual               2.01     1.418
## Number of obs: 38, groups:  sent_index, 19
##
## Fixed effects:
##             Estimate Std. Error      df t value Pr(>|t|)
## (Intercept)   2.0372     0.6405 18.0000   3.181  0.00518 **
## verb_type1    2.6455     0.2300 18.0000  11.503 9.94e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##            (Intr)
## verb_type1 0.000
```

- Statistically significant effect of the affirmativeness vs. negativeness of verbs.
- Greater effect than TOP.

## DAT

```
sub_data = subset(data_V_shika, shika_case == 'DAT')


m = lmer(
        surprisal_diff
            ~ verb_type
                + (1 | sent_index)
        ,
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: surprisal_diff ~ verb_type + (1 | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 128.9
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -1.1309 -0.7259  0.1149  0.5097  1.3964
##
## Random effects:
##  Groups      Name        Variance Std.Dev.
##  sent_index (Intercept) 4.020    2.005
##  Residual               1.256    1.121
## Number of obs: 32, groups:  sent_index, 16
##
## Fixed effects:
##              Estimate Std. Error      df t value Pr(>|t|)
## (Intercept)  0.08783    0.53897 15.00000   0.163    0.873
## verb_type1   2.67121    0.19811 15.00000  13.483 8.65e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##            (Intr)
## verb_type1 0.000
```

- Statistically significant effect of the affirmativeness vs. negativeness of verbs.
- Greater effect than TOP, similar to ACC.