# rnnpsycholing Japanese NPI (sentences with embedding and matrix shika)

*Takashi Morita*

## 目次

## Introduction

We are looking for a 2x2x2 interaction of:

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

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

e.g.

- TOP
  - 佐藤-{しか, は} 社長-が パーティ-に 友人-を {呼んだ, 呼ばなかった} と {思った, 思わなかった}。
  - Sato-{*shika*, TOP} CEO-NOM party-DAT friend-ACC {invited, didn't invite} that {thought, didn't think}.
- DAT
  - 同僚-{にしか, に} 佐藤は 社長-が パーティ-に 友人-を {呼んだ, 呼ばなかった} と {思った, 思わなかった}。

– colleague-{DAT-*shika*, DAT} Sato-TOP CEO-NOM party-DAT friend-ACC {invited, didn't invite} that {thought, didn't think}.

Why is this interesting?

1. A grammatical sentence with *shika* in the main clause must have a negative main verb.
   - A significant increase in surprisal of the affirmative main verbs must be predicted by the LSTM conditioned on the presence of *shika* if the learning is successful.
2. Negation of the embedded verb does not satisfy the *shika*'s grammatical condition.
   - No significant increase in surprisal of the affirmative embedded verbs given *shika* is expected for a successful learner.
   - Nor significant interaction between the main and embedded verbs given *shika* is expected for a successful learner.

## Load data

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


REGIONS = c('main_prefix', 'embedded_prefix', 'embedded_V', 'complementizer', 'main_V', 'end')


token_based_data_path = 'jp_shika_test_sentences_embedded_shika-in-main_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 = col_character(),
##   embed_V = col_character(),
##   main_V = 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, embed_V, main_V, shika_case) %>%
        summarise(surprisal=sum(surprisal)) %>%
        ungroup() %>%
    mutate(
        shika=factor(shika, levels=c("shika", "no-shika")),
        embed_V=factor(embed_V, levels=c("affirmative", "negative")),
        main_V=factor(main_V, levels=c("affirmative", "negative")),
        shika_case=factor(shika_case, levels=c("TOP", "DAT"))
        )

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

# Make sure that the dataframe is sorted appropriately.
# First by embed_V (affirmative vs. negative)
data_region_based = data_region_based[order(data_region_based$embed_V),]
# Then by main_V
data_region_based = data_region_based[order(data_region_based$main_V),]
# finally by sent_index
data_region_based = data_region_based[order(data_region_based$sent_index),]
```

## Embedded verb region

### Visualization

```r
# Focus on the V (verb) region.
data_V = subset(data_region_based, region == 'embedded_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 %>%
```
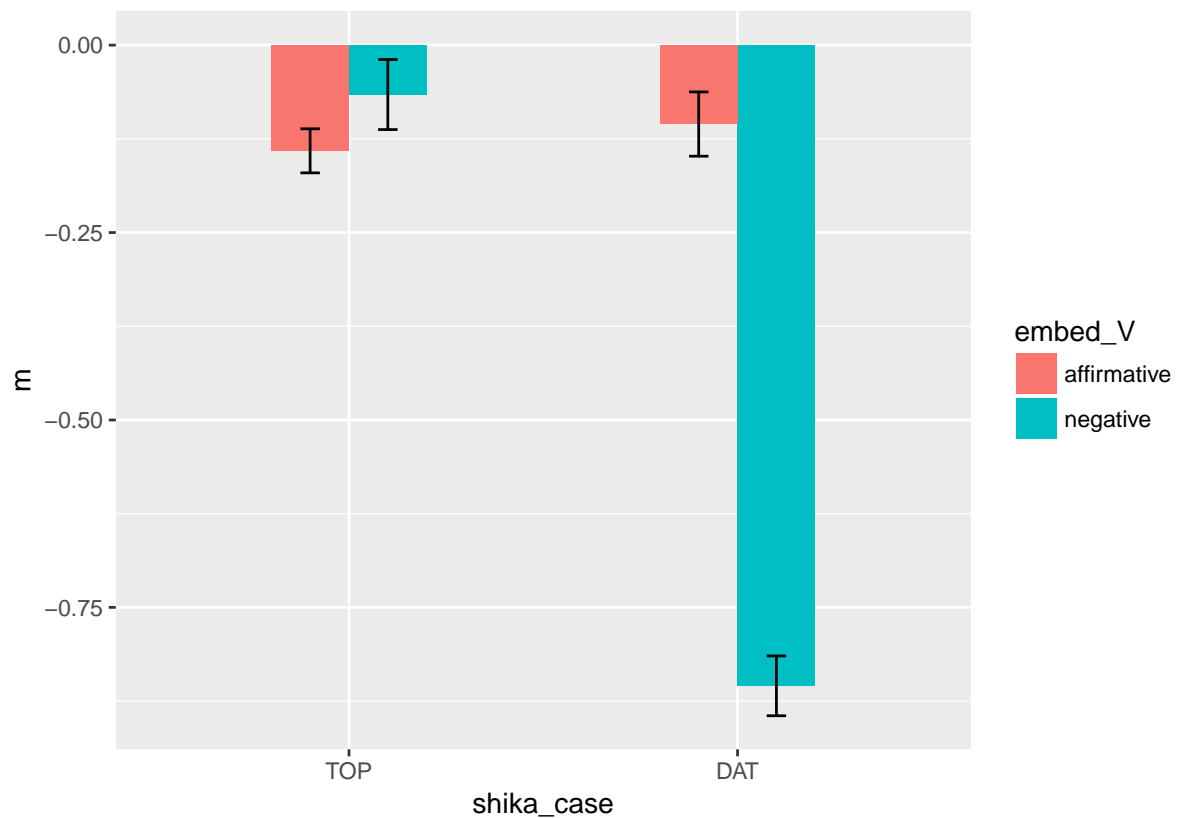
```
    group_by(embed_V, shika_case) %>%
    summarise(m=mean(surprisal_diff),
            s=std.error(surprisal_diff),
            upper=m + 1.96*s,
            lower=m - 1.96*s) %>%
    ungroup() %>%
    ggplot(aes(x=shika_case, y=m, ymin=lower, ymax=upper, width=0.4, fill=embed_V)) +
        geom_bar(stat = 'identity', position = "dodge") +
        geom_errorbar(position=position_dodge(0.4), width=.1)
```



## Regressions

### TOP

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


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

```
        )
summary(m)
```

## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: surprisal_diff ~ embed_V + (1 | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 3978.9
##
## Scaled residuals:
##      Min      1Q   Median      3Q      Max
## -2.59434 -0.57830  0.00579  0.58077  2.00527
##
## Random effects:
##  Groups     Name        Variance Std.Dev.
##  sent_index (Intercept) 0.3973   0.6303
##  Residual               0.1357   0.3683
## Number of obs: 2688, groups:  sent_index, 672
##
## Fixed effects:
##               Estimate Std. Error        df t value Pr(>|t|)
## (Intercept) -1.035e-01  2.533e-02  6.710e+02  -4.084 4.95e-05 ***
## embed_V1    -3.755e-02  7.104e-03  2.015e+03  -5.285 1.39e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##          (Intr)
## embed_V1 0.000

- Significant negative effect of embed_V (affirmativeness = 1).
    - Negative verbs cause more

DAT

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


m = lmer(
        surprisal_diff
            ~ embed_V
                + (1 | sent_index)
        ,
```

```
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula: surprisal_diff ~ embed_V + (1 | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 1811.6
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -3.5575 -0.3278  0.0332  0.4594  2.8311
##
## Random effects:
##  Groups     Name        Variance Std.Dev.
##  sent_index (Intercept) 0.2363   0.4861
##  Residual               0.1067   0.3266
## Number of obs: 1536, groups:  sent_index, 384
##
## Fixed effects:
##               Estimate Std. Error        df t value Pr(>|t|)
## (Intercept) -4.800e-01  2.617e-02 3.830e+02  -18.34   <2e-16 ***
## embed_V1     3.747e-01  8.334e-03 1.151e+03   44.96   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##          (Intr)
## embed_V1 0.000
```

- Significant positive effect of embed_V (affirmativeness = 1).

## Main verb region

### Visualization

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


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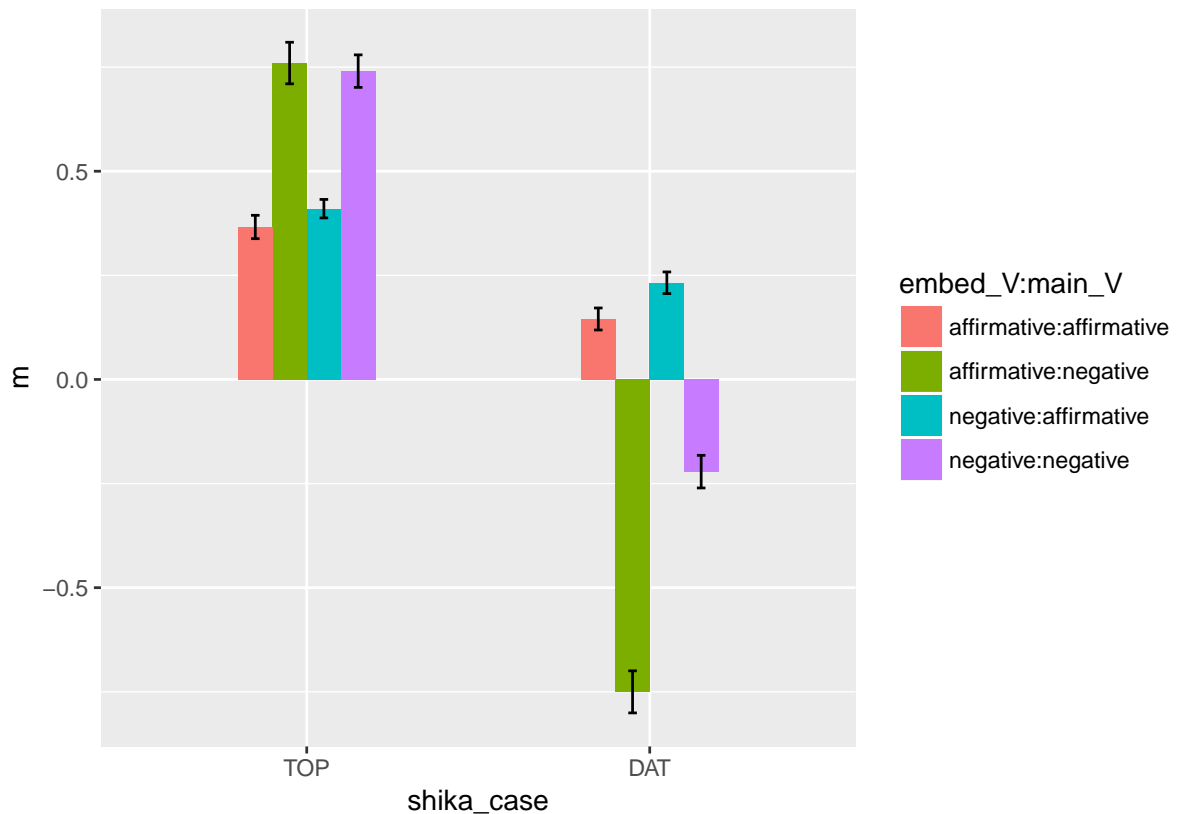
```
data_V_no_shika = subset(data_V, shika == 'no-shika')
data_V_shika$surprisal_diff = data_V_shika$surprisal - data_V_no_shika$surprisal



# Visualize the difference in surprisal increase/dicrease between affirmative vs. negative verbs.
data_V_shika %>%
    group_by(embed_V, main_V, 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=embed_V:main_V)) +
        geom_bar(stat = 'identity', position = "dodge") +
        geom_errorbar(position=position_dodge(0.4), width=.1)
```



- TOP
    - Increas in surprisal in every condition.
- DAT
    - Small but expected signs of changes.

Regressions

TOP

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


m = lmer(
        surprisal_diff
            ~ embed_V * main_V
                + (embed_V + main_V | sent_index)
        ,
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula:
## surprisal_diff ~ embed_V * main_V + (embed_V + main_V | sent_index)
##     Data: sub_data
##
## REML criterion at convergence: 1604.5
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -4.9448 -0.3444 -0.0138  0.3504  4.8049
##
## Random effects:
##  Groups     Name        Variance Std.Dev. Corr
##  sent_index (Intercept) 0.13922  0.3731
##             embed_V1    0.01594  0.1262    0.52
##             main_V1     0.05716  0.2391   -0.59 -0.18
##  Residual               0.01922  0.1386
## Number of obs: 2688, groups:  sent_index, 672
##
## Fixed effects:
##                   Estimate Std. Error         df t value Pr(>|t|)
## (Intercept)       0.569087   0.014640 671.000343  38.873  < 2e-16 ***
## embed_V1         -0.006159   0.005556 671.000200  -1.109    0.268
## main_V1          -0.180980   0.009603 670.999659 -18.846  < 2e-16 ***
## embed_V1:main_V1 -0.015856   0.002674 670.999947  -5.929 4.87e-09 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##             (Intr) emb_V1 man_V1
## embed_V1     0.451
## main_V1     -0.559 -0.156
## embd_V1:_V1  0.000  0.000  0.000
```

- No significant effect of embed_V (affirmativeness = 1).
- Significant negative effect of main_V (affirmativeness = 1).
- Significant negative interaction.

DAT

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


m = lmer(
        surprisal_diff
            ~ embed_V * main_V
                + (embed_V + main_V | sent_index)

        ,
        data=sub_data
        )
summary(m)
```

```
## Linear mixed model fit by REML. t-tests use Satterthwaite's method [
## lmerModLmerTest]
## Formula:
## surprisal_diff ~ embed_V * main_V + (embed_V + main_V | sent_index)
##    Data: sub_data
##
## REML criterion at convergence: 600.4
##
## Scaled residuals:
##     Min      1Q  Median      3Q     Max
## -4.5474 -0.3487 -0.0140  0.3562  3.0530
##
## Random effects:
##  Groups     Name        Variance Std.Dev. Corr
##  sent_index (Intercept) 0.068481 0.26169
##             embed_V1    0.008776 0.09368   0.25
##             main_V1     0.032726 0.18090  -0.56 -0.49
##  Residual               0.026540 0.16291
## Number of obs: 1536, groups:  sent_index, 384
```

9

```
##
## Fixed effects:
##                   Estimate Std. Error        df t value Pr(>|t|)
## (Intercept)       -0.148561   0.013986 382.999963  -10.62   <2e-16 ***
## embed_V1          -0.153967   0.006335 383.000063  -24.30   <2e-16 ***
## main_V1            0.337250   0.010124 383.000058   33.31   <2e-16 ***
## embed_V1:main_V1   0.110475   0.004157 382.999925   26.58   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##           (Intr) emb_V1 man_V1
## embed_V1   0.182
## main_V1   -0.489 -0.337
## embd_V1:_V1  0.000  0.000  0.000
```

- Significant negative effect of embed_V (affirmativeness = 1).
- Significant positive effect of main_V (affirmativeness = 1).
- Significant positive interaction.