

# Project 2

Xiaonan Hu, Lee Stovall, Liqi Zhu

---

## Part 1 ASJP Dataset

### 1. Introduction

The purpose of Part 1 is to find the overall status for those languages we choose based on the population that use the language. ([http://www.vistawide.com/languages/language\\_families\\_statistics1.htm](http://www.vistawide.com/languages/language_families_statistics1.htm))

Language Family	Speaker Percentage
Indo-European	44.78%
Sino-Tibetan	22.28%
Niger-Congo	6.26%
Afro-Asiatic	5.93%
Austronesian	5.45%
Dravidian	3.87%
Japanese	2.16%
Austro-Asiatic	1.77%
Tai-Kadai	1.37%
<b>Total</b>	<b>93.87%</b>

We managed to download the dataset for each language family from ASJP database.

### 2. Method & Processing Data

Our thought is to compare the Distance data within each language family. One dataset is with the numbers and

the other is not.

```
1. ##### Cygwin #####
2. egrep -v '\sone|\stwo' < Indo-European.txt > IEwonum.txt
3. egrep -v '\sone|\stwo' < Austro-Asiatic.txt > AuAwonum.txt
4. egrep -v '\sone|\stwo' < Austronesian.txt > Awonum.txt
5. egrep -v '\sone|\stwo' < Dravidian.txt > Dwounum.txt
6. egrep -v '\sone|\stwo' < Japanese.txt > Jwonum.txt
7. egrep -v '\sone|\stwo' < Niger-Congo.txt > NCwonum.txt
8. egrep -v '\sone|\stwo' < Sino-Tibetan.txt > STwnum.txt
9. egrep -v '\sone|\stwo' < Tai-Kadai.txt > TKwonum.txt
10. egrep -v '\sone|\stwo' < Afro-Asiatic.txt > AAwonum.txt
11. # Exclude numbers in dataset > Datasets without numbers.
```

With Programs for calculating ASJP distance matrices (<http://asjp.clld.org/software>). We managed to get the distance matrix for each language family.

```
1. ##### Power Shell #####
2. asjp62 < listss17.txt > output1.txt
3. asjp62 < Indo-European.txt > IEall.txt
4. asjp62 < Austro-Asiatic.txt > AuAall.txt
5. asjp62 < Austronesian.txt > Aall.txt
6. asjp62 < Dravidian.txt > Dall.txt
7. asjp62 < Japanese.txt > Jall.txt
8. asjp62 < Niger-Congo.txt > NCall.txt
9. asjp62 < Tai-Kadai.txt > TKall.txt
10. asjp62 < Sino-Tibetan.txt > STall.txt
11. asjp62 < Afro-Asiatic.txt > AAall.txt
12. # Overall view of whole dataset and different language families with n
    umbers.
13. asjp62 < IEwonum.txt > IEall_won.txt
14. asjp62 < AuAwonum.txt > AuAall_won.txt
15. asjp62 < Awonum.txt > Aall_won.txt
16. asjp62 < Dwounum.txt > Dall_won.txt
17. asjp62 < Jwonum.txt > Jall_won.txt
18. asjp62 < NCwonum.txt > NCall_won.txt
19. asjp62 < TKwonum.txt > TKall_won.txt
20. asjp62 < STwnum.txt > STall_won.txt
21. asjp62 < AAwonum.txt > AAall_won.txt
22. # Overall view of whole dataset and different language families without
    numbers.
```

```
1. cat IEall*.txt > IE.txt
```

```
2. cat AuAall*.txt > AuA.txt
3. cat Aall*.txt > A.txt
4. cat Dall*.txt > D.txt
5. cat Jall*.txt > J.txt
6. cat NCall*.txt > NC.txt
7. cat TKall*.txt > TK.txt
8. cat STall*.txt > ST.txt
9. cat AAall*.txt > AA.txt
10.
11. sed 's/ \+/,/g' IE.txt > IE.csv
12. sed 's/ \+/,/g' AuA.txt > AuA.csv
13. sed 's/ \+/,/g' A.txt > A.csv
14. sed 's/ \+/,/g' D.txt > D.csv
15. sed 's/ \+/,/g' J.txt > J.csv
16. sed 's/ \+/,/g' NC.txt > NC.csv
17. sed 's/ \+/,/g' TK.txt > TK.csv
18. sed 's/ \+/,/g' ST.txt > ST.csv
19. sed 's/ \+/,/g' AA.txt > AA.csv
20. # Arrange data into CSV file
```

With the combination and calculation of the .CSV files, we are able to compare the ASJP Distance of each language family with and without numbers.

## 3. Results

Example of Japanese

As showed above, values are mesured ASJP Distance data. We calculated the average standard for each family language family and then compare the Distance matrix with numbers and the one without.

Then we repeat the process for the other language families.

The results of all language families are followed:

	Speaker Percentage	ASJP Distance With Numbers	ASJP Distance Without Numbers	Difference
Indo-European	44.78%	85.56	85.84	0.280495
Sino-Tibetan	22.28%	88.27	88.46	0.191502
Niger-Congo	6.26%	91.43	91.32	-0.1074
Afro-Asiatic	5.93%	91.32	91.16	-0.15255
Austronesian	5.45%	85.31	85.40	0.089109
Dravidian	3.87%	85.31	57.28	-28.0276
Japanese	2.16%	48.82	48.50	-0.32087
Austro-Asiatic	1.77%	88.80	81.34	-7.4635
Tai-Kadai	1.37%	75.99	75.61	-0.37409
Average	10.43%	82.31	78.32	-3.99
Total	93.87%	740.80	704.92	-35.8849

As shown above, the first thing we discovered is that language families have different features. Japanese was obviously the language family containing the most similar languages since ASJP Distance is the smallest. While all the other language families show different level of variance.

About the influence of numbers, Dravidian's Distance drops significantly, which means the number changes a lot in Dravidian Languages. Thus to say all the language family with the positive numbers of Difference means number changes less than other words. And Negative number indicates number changes more.

As a conclusion, in Indo-European, Sina-Tibetan, Austronesian language families, numbers changes less. While in other language families numbers may change more than average.