Linguistic Data Analysis - Final Project

Benedict Wuethrich

2024-07-10

Introduction

On the Australian continent there are over 333 reported languages which can be roughly categorized into either Pama-Nyungan languages or Non-Pama-Nyungan languages. It is important to note that the latter does not imply any genealogical link of the included languages, while the former has been shown as a cohesive language family. Non-Pama-Nyungan is more of a collective term for Australian languages that aren't thought of as Pama-Nyungan and actually contain 27 different language families. Exactly which languages belong to which category is an ongoing debate.

Pama-Nyungen languages have been spoken for over 5000 years, make up over 306 different identified languages and the speakers cover about 80% of the landmass (Bouckaert 2018, 741). Given this long history, it is likely that the Pama-Nyungan languages have been competing against Non-Pama-Nyungan languages for a long time and ended up limiting the Non-Pama-Nyungan languages to the northernmost part of the continent. Nowadays, both are having to compete against the English language.

One of the main actors in describing Australian languages is Robert M. W. Dixon. With "The Languages of Australia", he has published a comprehensive description of Australian languages as early as 1980 and another study of Australian Languages in 2004 in his monograph "Australian Languages. Their Nature and Development". But these works were not without their controversy. Dixon vehemently denied the applicability of the comparative method on Australian Languages, which lead to a rather large amount of discussion surrounding his works. Claire Bowern and Harold Koch wrote "Australian Languages. Classification and the comparative method". Dixon's skepticism is described as an erroneous phylogenetic assessment which is "so bizarrely faulted, and such an insult to the eminently successful practitioners of Comparative Method Linguistics in Australia, that it positively demands a decisive riposte." (O'Grady and Hale in Bowern and Koch, 2004: 69), which serves as a perfect summary for what the book tries to achieve.

In more recent publications, like the paper by Bouckaert et al. (2018), computer-assisted methods have gained ground. This paper specifically deals with the genealogy of Pama-Nyungan languages. Using their data consisting of basic vocabulary of 306 different Pama-Nyungan languages, they created a phylogenetic tree that shows the diversification of the Pama-Nyungan language family. Additionally, the same data was used to find a likely homeland of the family, which they postulated to be in the Gulf Plains region.

In this paper I will be looking at the geographical distribution of the Australian languages and how they compare to each other in terms of categorical features such as word order, use of reduplication and more. To make this both simpler to understand as well as simpler to execute, I will be adopting the reductionist point of view of pooling the Non-Pama-Nyungan languages into a single category. Furthermore, for certain features, I will focus solely on the Pama-Nyungan languages.

Data

The data used for this project originates from the cldf data by WALS. I initially intended to use the phoible cldf data as well, but ended up deciding against it as it did not add an appropriate amount of value to the final product.

The data wrangling for this project was manageable, but took a lot of work to get started. The biggest challenge posed was figuring out how to load cldf dataset into R and how to manipulate it appropriately. I struggled with that part for a while as it also had me refresh my memory on how to manipulate dataframes eficiently and how to structure the creation of new dataframes from the given sets without creating too much redundancy. I created a seperate Rscript file that entails all the necessary data wrangling and manipulation for the different plots and maps used in this project. The other scripts rely on this Data Wrangling script to be run first, so as to load the required dataframes. This Rmarkdown file compiles them all into one coherent document, though without saving the plots and maps again. The seperate scripts can be found in the corresponding GitHub repository, a link to which can be found at the end of this document.

I started out by loading in the cldf using the rcldf package. From this a wide dataframe was created using the Value Table and filtering for Australian languages only.

```
##loading in data from phoible and WALS and limiting them to languages spoken in Australia
#Let's start with WALS
wals_cldf <- cldf(here("data/WALS/cldf"))</pre>
#this loads the cldf data, but is not yet manipulable. Do to this, we need to get it into a dataframe o
summary(wals cldf)
## A Cross-Linguistic Data Format (CLDF) dataset:
## Name: The World Atlas of Language Structures Online
## Path: C:/Users/Benedict/Documents/GitHub/LDAR-Project_BW/LDAR-FP-BW/data/WALS/cldf
## Type: http://cldf.clld.org/v1.0/terms.rdf#StructureDataset
## Tables:
##
     1/12: areas.csv (3 columns, 11 rows)
     2/12: CodeTable (6 columns, 1143 rows)
##
##
     3/12: ContributionTable (11 columns, 152 rows)
##
     4/12: contributors.csv (4 columns, 55 rows)
     5/12: countries.csv (2 columns, 192 rows)
##
     6/12: ExampleTable (8 columns, 3907 rows)
##
     7/12: language_names.csv (4 columns, 7377 rows)
##
     8/12: LanguageTable (17 columns, 3573 rows)
##
##
     9/12: MediaTable (8 columns, 153 rows)
     10/12: ParameterTable (5 columns, 192 rows)
##
     11/12: TreeTable (8 columns, 254 rows)
##
     12/12: ValueTable (8 columns, 76475 rows)
##
## Sources: 7373
#This creates a tibble that will be manipulated further at different points.
wals_value <- as.cldf.wide(wals_cldf, "ValueTable") %>%
 filter(Macroarea == "Australia")
## Joining Parameter_ID -> ParameterTable -> ID
## Joining Code_ID -> CodeTable -> ID
## Joining Language_ID -> LanguageTable -> ID
## Joining Example_ID -> ExampleTable -> ID
```

glimpse(wals_value)

```
## Rows: 5,096
## Columns: 40
## $ ID
                         <chr> "26A-abn", "30A-abn", "31A-abn", "32A-abn",~
                         <chr> "abn", "abn", "abn", "abn", "abn", "abn", "~
## $ Language ID. ValueTable
                         <chr> "26A", "30A", "31A", "32A", "33A", "35A", "~
## $ Parameter_ID.ValueTable
                         <chr> "2", "1", "1", "1", "2", "4", "2", "3", "1"~
## $ Value
## $ Code_ID
                         <chr> "26A-2", "30A-1", "31A-1", "32A-1", "33A-2"~
## $ Comment.ValueTable
                         <chr> NA, NA, NA, NA, NA, NA, NA, "<div class=\"E~
## $ Source.ValueTable
                         <chr> "Hercus-1994[passim]", "Hercus-1994", "Herc~
## $ Example ID
                         ## $ Name.ParameterTable
                         <chr> "Prefixing vs. Suffixing in Inflectional Mo~
## $ ColumnSpec
                         <chr> "26", "30", "31", "32", "33", "35", "39", "~
## $ Chapter_ID
                         <chr> "26A", "30A", "31A", "32A", "33A", "35A", "~
## $ Parameter_ID.CodeTable
                         <chr> "Strongly suffixing", "None", "No gender", ~
## $ Name.CodeTable
                         <chr> "Predominantly suffixing", "None", "No gend~
## $ Description.CodeTable
## $ Number
                         <int> 2, 1, 1, 1, 2, 4, 2, 3, 1, 6, 1, 1, 1, 2, 4~
## $ icon
                         <chr> "c0000dd", "cfffffff", "cfffffff", "cfffffff", "
## $ Name.LanguageTable
                         <chr> "Arabana", "Arabana", "Arabana", "Arabana",~
                         <chr> "Australia", "Australia", "Australia", "Aus~
## $ Macroarea
## $ Latitude
                         <dbl> -28.25, -28.25, -28.25, -28.25, -28.25, -28.
## $ Longitude
                         <dbl> 136.25, 136.25, 136.25, 136.25, 136.25, 136.
## $ Glottocode
                         <chr> "arab1267", "arab1267", "arab1267", "arab12~
                         <chr> "ard", "ard", "ard", "ard", "ard", "ard", "~
## $ ISO639P3code
                         <chr> "Pama-Nyungan", "Pama-Nyungan", "Pama-Nyung~
## $ Family
## $ Subfamily
                         <chr> "Central Pama-Nyungan", "Central Pama-Nyung~
## $ Genus
## $ GenusIcon
                         ## $ ISO_codes
                         <chr> "ard", "ard", "ard", "ard", "ard", "ard", "~
## $ Samples_100
                         <lg1> FALSE, FALSE, FALSE, FALSE, FALSE, F~
                         <lg1> FALSE, FALSE, FALSE, FALSE, FALSE, F~
## $ Samples_200
                         <chr> "AU", "AU", "AU", "AU", "AU", "AU", "AU", "AU", "~
## $ Country_ID
## $ Source.LanguageTable
                         <chr> "Hercus-1994", "Hercus-1994", "Hercus-1994"~
## $ Parent_ID
                         <chr> "genus-centralpamanyungan", "genus-centralp~
                         ## $ Language_ID.ExampleTable
## $ Primary_Text
                         ## $ Analyzed_Word
                         ## $ Gloss
                         ## $ Translated Text
                         ## $ Meta_Language_ID
                         ## $ Comment.ExampleTable
```

Next up was creating seperate dataframes that specifically only included either Pama-Nyungan or Non-Pama-Nyungan languages. Additionally, the relevent column was renamed to be more descriptive and therefore easier to work with.

```
##Word Order Data Wrangling
#creates two seperate df, for PN and nPN respectively.
wals_worder_PN <- wals_value %>%
```

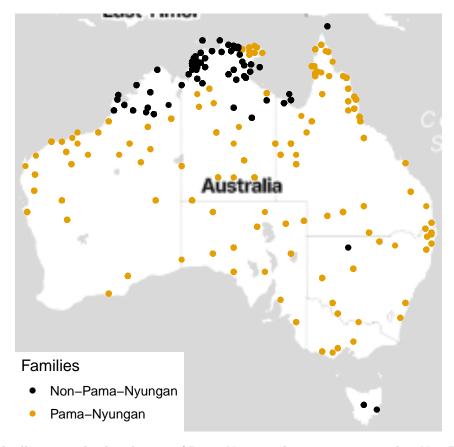
Analysis

First off, let's see the distribution of Pama-Nyungan languages compared to Non-Pama Nyungan languages. To do this, I first load a map of Australia through stadiamap. Next, I modify an existing dataframe in order to make the legend of the final map more easily readable. Using this new dataframe, I plot the locations of the languages onto the map.

i © Stadia Maps © Stamen Design © OpenMapTiles © OpenStreetMap contributors.

```
ggmap(map_AUS)
```



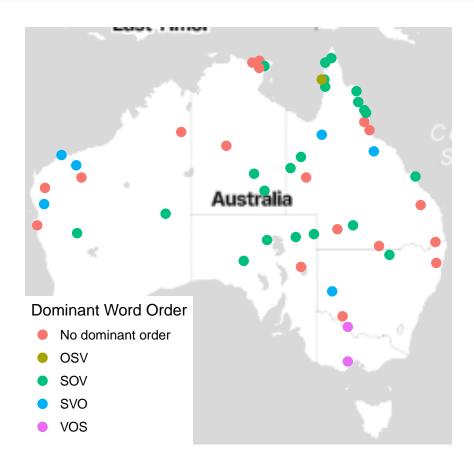


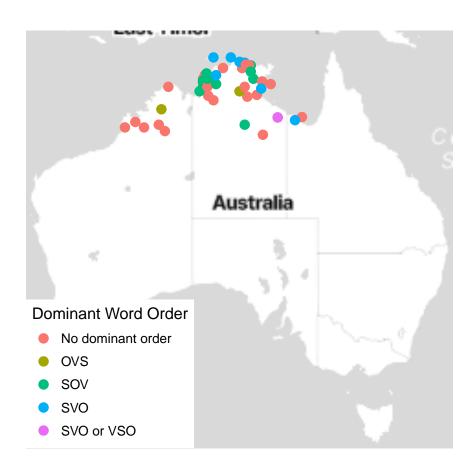
This map clearly illustrates the distribution of Pama-Nyungan languages compared to Non-Pama-Nyungan languages. The latter is, with a few exceptions, visibly confined to the northernmost part of Australia, whereas the Pama-Nyungan languages are spread all throughout the continent.

Word Order

The word order of different languages is described as "perhaps the single most frequently cited typological feature of languages" by Dryer and Haspelmath on the website concerning the Chapter about word order. I do not intend on making an exception in this matter either. To gain insight on the matter for the context of the Australian languages, I will be creating seperate maps for both Pama-Nyungan and Non-Pama-Nyungan, as well as seperate bar plots to explore the frequency of each within the respective categories.

```
labs(color = "Dominant Word Order")
map_PN_worder
```



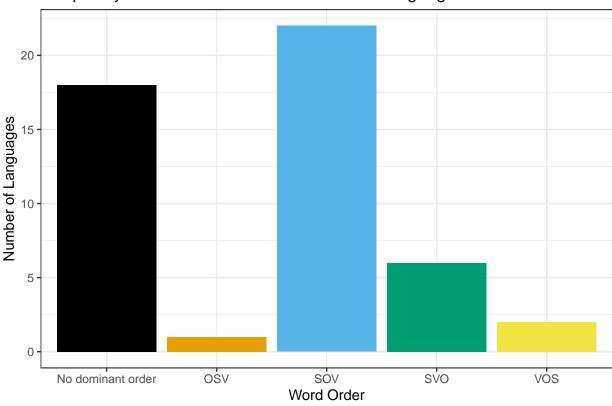


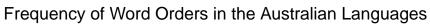
 ${\tt \#Note: the \ size \ of \ the \ dots \ was \ reduced \ compared \ to \ the \ code \ in \ Maps. R, \ so \ as \ to \ make \ the \ individual \ dots \ dots$

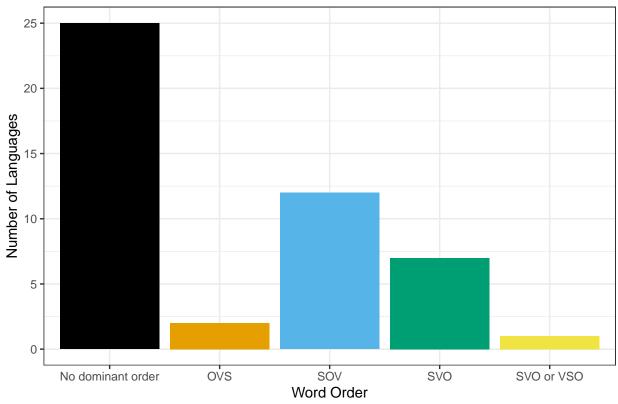
```
#PLOTS
##Word Order Plots
###Word Order PN
histo_PN_worder <- ggplot(data = wals_worder_PN,
                               aes(x = WordOrder,
                                   fill = WordOrder)) +
  geom_histogram(stat = "count") + # defining the type of plot
  labs(y = "Number of Languages",
       x= "Word Order") + # renaming the axes
  theme_bw() + # adjusting the theme
  scale_fill_colorblind(guide = FALSE) +
  ggtitle("Frequency of Word Orders in the Australian Languages") # adding a descriptive title
## Warning in geom_histogram(stat = "count"): Ignoring unknown parameters:
## `binwidth`, `bins`, and `pad`
histo_PN_worder #getting a preview
## Warning: The `guide` argument in `scale_*()` cannot be `FALSE`. This was deprecated in
```

```
## ggplot2 3.3.4.
## i Please use "none" instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.
```

Frequency of Word Orders in the Australian Languages







Bibliography

Bouckaert, Remco R., Claire Bowern & Quentin D. Atkinson (2018). The origin and expansion of Pama–Nyungan languages across Australia. Nature Ecology & Evolution (2), 741–749.

Bowern, Claire / Koch, Harold (Hrsg.) (2004): Australian Languages. Classification and the comparative method. Amsterdam / Philadelphia: John Benjamins.

Dixon, R.M.W. (1980): The Languages of Australia. Cambridge: Cambridge University Press.

Dixon, R.M.W. (2004): Australian Languages. Their Nature and Development. Cambridge: Cambridge University Press.

Greenhill, Simon (2024). rcldf: rcldf - Read Linguistic Data In The Cross Linguistic Data Format (CLDF)_. R package version 1.2.0, commit 3979a89dbe4db653873caca62212fc07ebb966e9, https://github.com/SimonGreenhill/rcldf

Kahle, D., Wickham, H. ggmap: Spatial Visualization with ggplot2. The R Journal, 5(1), 144-161. http://journal.r-project.org/archive/2013-1/kahle-wickham.pdf

Matthew S. Dryer. (2013) Order of Subject, Object and Verb. In: Dryer, Matthew S. & Haspelmath, Martin (eds.) WALS Online (v2020.3) [Data set]. Zenodo. https://doi.org/10.5281/zenodo.7385533 (Available online at http://wals.info/chapter/81, Accessed on 2024-04-16.)

Moran, Steven & McCloy, Daniel (eds.) 2019. PHOIBLE. Jena: Max Planck Institute for the Science of Human History. (Available online at https://phoible.org)

O'Grady, G. N. (1998). Toward a Proto-Pama-Nyungan Stem List, Part I: Sets J1-J25. Oceanic Linguistics, 37(2), 209–233.

Schmidt, W. (1919). Die Gliederung der australischen Sprachen: geographische, bibliographische, linguistische Grundzüge der Erforschung der australischen Sprachen. Mechitharisten-Buchdruckerei.

Wickham, H et al. (2019). Welcome to the tidyverse. Journal of Open Source Software, 4 (43), 1686. doi:10.21105/joss.01686 https://doi.org/10.21105/joss.01686

Zuckermann, G. et al. (2021) LARA in the Service of Revivalistics and Documentary Linguistics: Community Engagement and Endangered Languages. Proceedings of the Workshop on Computational Methods for Endangered Languages (1), 13-23.

Everything related to this project can be found in this GitHub repository.