

NLP Assisted Analysis of Folk Taxonomy

An examination of the Matukar language

Jonathan Gluck

Abstract

Folk taxonomies are powerful cultural tools for the categorization and utilization of the world in which a people live. The English language, for example, has a few folk taxa remaining; including *pets*, *farm animals*, and *evergreens*. Folk taxa are categories or logical groupings, usually referring to nature, which may have social and cultural relevance, but not necessarily possessing any scientific relatedness amongst their members. They are useful in day-to-day dealings with the environment, providing a catalogue grouped by salient features. Finding a language's folk taxonomy can often be difficult, with the lines drawn between categories often not readily apparent. With this work I examine the theory behind folk taxonomic classification and attempt to devise methods for unearthing folk taxonomies with the help of Natural Language Processing.

The subject language of this inquiry is Matukar. Matukar is an Austronesian language of Papua New Guinea, spoken by only about 430 villagers on the North Eastern coast. The language is spoken in a rural area of costal Papua New Guinea and contains many of the ambient sounds of their surroundings. It is a language threatened by the rising popularities of English and the local Creole ,Tok Pisin. The

folk taxonomy of Matukar has never before been examined, and is the focus of this work.

The job of unearthing a folk taxonomy involves sifting through large numbers of dictionary entries and searching for patterns and similarities in word form, be they morphological or phonetic. Procedures, like these, which make use of large amounts of data are perfectly suited to Natural Language Processing, or NLP for short. NLP is the subfield of Computer Science most concerned with language and its use. With the help of NLP it is possible to process quantities of data that might otherwise be prohibitive for hand analysis.

Often members in a folk taxon have similar names, or exhibit internal patterns. (Berlin, Breedlove, & Raven, p. 216) One such example is the use of *fish* to group marine life in *jellyfish* and *goldfish*. In order to find such examples I use the NLP tool of string similarity. This involves comparing the distance between any two words' similarities and selecting for those that pass a certain threshold. This tool should provide a list of similar words in a target language, revealing similar folk taxa.

While members in a given folk taxonomy may not directly map to English's science influenced taxonomy, many of the borders between folk taxa are influenced by their members' higher level categories. (Hunn, pp. 830-831) Imposing English's taxonomy onto a target language might provide helpful categories within which to look for morphological similarities. In order to do this, I implement automatic semantic tagging using WordNet in concert with the English gloss for each Matukar word.

With the assistance of NLP the examination of folk taxonomies may be streamlined, providing linguists with a starting point with which to theorize folk taxa. I show the results of these tools on the Matukar Language.

Introduction

The range of human interaction, both in natural and social spheres, is vast. Even so, we humans are able to wrap our minds around the complex world in which we survive. The catalogue of discrete objects maintained in the human mind is of astounding length, so much so that the mere listing of a subset of this catalogue, for example names of familiar games, is rendered impossible. Access to this entire list at once is not possible. Yet, if “*Hop Scotch*” or “*Mother May I*” are referenced, the audience, so long as it has met with these games before, knows immediately not only that they are games, but also the environment in which they might be played and a myriad of other details. Accessing this knowledge is possible because of the human process of categorization. Humans observe the dynamics of their surroundings and file away their daily experiences for later use.

One specific, useful type of categorization is the Folk Taxonomy, or Folksonomy for short. Folk taxonomies are cultural methods developed over time for the classification and compartmentalization, of the day-to-day experiences of human life. They are traditionally biological, although there should be no reason for folk taxonomies to be confined to biology only.¹ They allow an understanding of species and how they relate to one another. They are culturally relevant tools, and

though they are not necessarily standard throughout a culture, they are a powerful tool to allow for the organization and control of the surrounding environment.

The goal of this project is to examine the theory behind folk taxonomies, and then analyze one language, Matukar, for clues pointing to possible folksonomies. The search for folksonomies will be undertaken with the help of Natural Language Processing tools operating on the Matukar Online Talking Dictionary.

A Survey of Matukar

Matukar is an endangered language of Papua New Guinea, spoken in two villages in the Madang Provinceⁱⁱ. The language, at current count, has about 430 speakers, including both “experienced elders and children.” (Harrison, Anderson, & Mathieu-Reeves, 2010) Matukar is endangered, as a language, because of the continual rising popularity of English and of Tok Pisin, the local creole and most common language of Papua New Guinea. (The Central Intelligence Agency, 2009) While there is much that is not known about the language, we do have some pertinent facts which may impact its potential folksonomy. Matukar’s villages are situated along the coast line; thus common animal categories and species might range from aquatic to terrestrial to avian in form. An interesting feature of the language is that it contains many onomatopoeic words for living things. (Harrison, Anderson, & Mathieu-Reeves, 2010) It is also important to note that the main agricultural products of the area are: palm, sweet potatoes, shellfish, poultry, and pork. (The Central Intelligence Agency, 2009) These products bear keeping in mind

as we undertake analysis of the language. The more culturally relevant a word, the more likely it is to have an instance of taxonomic import.

The medium through which I will explore the Matukar language is The Matukar Online Talking Dictionary. This is a dictionary of some 3,045 entries with associated audio recordings. There are no other published corpora of Matukar. It should be noted that this is not a large dictionary and it was not created with the goal of folk biological elicitation in mind, so results are likely to be incomplete.

Three Theories of Folk Taxonomy

The importance of human classification has engendered much debate. How does the human mind structure information? How does this information relate to the concrete biological hierarchy of modern scientific taxonomy? With what mindset should folk taxonomies be approached? In this section I will examine the arguments of three scholars on these issues and present their proposed folk taxonomic models.

Extendable Hierarchical Model

Brent Berlin is an American anthropologist most famous for his work on color terms. Berlin outlines a number of points on the subject of folk taxonomies. It is his belief that the similarities between folk taxonomies and scientific taxonomies have been ignored, and that this should change. Berlin begins by stating, "In all languages it is possible to isolate groupings of organisms known as 'taxa'" (Berlin, Breedlove, & Raven, 1973, p. 214) These taxa are grouped into small, ethno-biological categories, which are arranged into a hierarchy. These taxonomic

categories are as follows: unique beginner, life form, generic, specific, and varietal. Taxa of the same category tend to occur at the same level, but this is not required. They are diagrammed below with examples for each category in Figure 1.

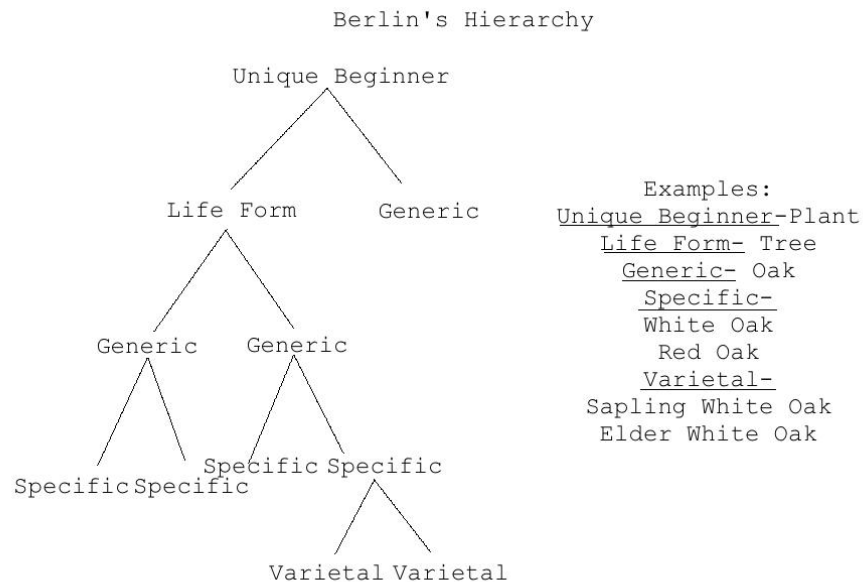


Figure 1: Berlin's Model

According to Berlin, the unique beginner category often goes unnamed in folk taxonomies. This unique beginner is something like “organism,” “animal,” or “plant.” Directly underneath the unique beginner are the life forms. Life forms tend to be few but important. Most taxa fit into one of the life forms. Berlin states, of generics, that they are more numerous than any other taxon. Most generics are immediately included as a child of some life form. Generics are the most important taxa for daily life. They are the taxa that are most quickly acquired by children. Sometimes generics are found without a parent life form class. In these cases, the generic is usually a borrowed word. (Berlin, Breedlove, & Raven, p. 220)

Once Berlin lays out his taxonomic hierarchy, he undertakes a short explanation of the formation of these words. He shows that, in his system, all taxa, with the exceptions of specific and varietal, are denoted by “primary lexemes.” Specific and varietal taxa are denoted by “secondary lexemes.” (Berlin, Breedlove, & Raven, 1973, p. 216) Primary lexemes tend to be single words and can be either analyzable (blueberry) or un-analyzable (spruce.) Secondary lexemes tend to be made up of two words, a descriptive word and a primary lexeme from another taxa, for example “blue spruce.”

Berlin’s arguments are compelling. The true utility of his hierarchy stems from its flexibility. He attempts, through his arguments, to find a model that is a compromise of several older models. In doing so he creates a truly extensible system.

Central Decentralized Model

Eugene Hunn is an American anthropologist who has a special focus on the cognitive aspects of ethno-biology. He is of the opinion that ethno-biology as a field has lost sight of the importance of examining the utility of folk taxonomies. He exhibits a strong belief that folk taxonomies are products of necessity and thus intrinsically utilitarian. In this vein, he gives a nod to Berlin who acknowledges that folksonomies are often affected by “cultural significance” (Berlin, Breedlove, & Raven, 1973, p. 839) Hunn explains that one reason for the utilitarian basis of folksonomies is that there is an information processing limitation that is imposed by the sheer number of possible items to classify. Thus, we must process those species that are the most useful first.

Hunn forgoes the hierarchical model for a centralized/decentralized model. He explains that the central categories are the easiest to recall. They are polythetic, determined by several optional characteristics. Non-central categories are both artificial and monothetic; members of these sets must subscribe to strict properties. This system is diagrammed, with examples, in figure 2.

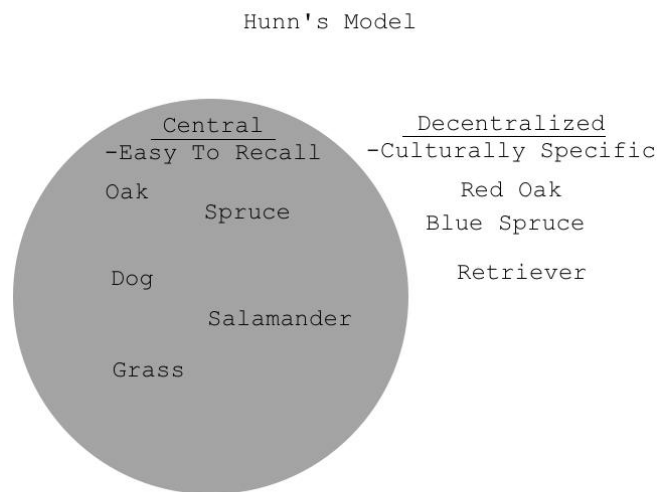


Figure 2: Hunn's Model

Hunn believes that Berlin might be attempting to jam these "central" categories into his generic taxa. This, to Hunn, seems "awkward" (Hunn, 1982, p. 836), as the generic class, in Berlin's hierarchy, is often found at several different locations, superordinate and subordinate to the generic taxa level. Hunn also highlights an issue with Berlin's parallels between scientific and folk hierarchy, that the folk taxa, "bird," might be entirely different from the scientific taxa of the same name. The folk taxa, for example, might refer to "environmental or aerial habitats,"

(Hunn, 1982, p. 838) while the scientific taxa are concerned with biological relatedness.

Hunn's central/decentralized model is an appealing alternative to Berlin's hierarchy. Hunn is concerned by the overwhelming focus on folk taxonomies as examples of "classification for its own sake." (Hunn, 1982, p. 831) Hunn proposes that the utility of each word in a given taxonomy be examined closely before attempts are made at compiling a model of that folk taxonomy.

Concrete Hierarchical Model

Scott Atran is a French American anthropologist. He is concerned with universal concepts in human thought and society. He currently studies biological classification in the mind. Atran believes that the system of classification present in folk taxonomies is "a cognitive mapping that places living-kind categories in a structure of absolute levels, which may... correspond to different levels of reality." (Atran, 1995, p. 141) Based on this statement, Atran's theory is more akin to Berlin's hierarchy than to Hunn's central/decentralized model. Additionally, it suggests that Atran believes folksonomies have a basis on reality. Atran states that the concept of folk taxonomies is hinged on the belief that variation not only exists in nature, but that it divides down salient lines. (Atran, 1995, p. 135) Humans develop taxonomic classes and imbue them with qualities learned from "naturalness." (Atran, 1995, p. 137) Naturalness, in this case, refers to the quality of an object, which belongs to a category, being associated with the rules governed by that category. (e.g.: even a pygmy elephant is cognized as a huge animal by being an elephant.) Atran points out that folk biological taxonomies are special in that they

have this quality of naturalness. Taxonomies of artifice do not exemplify this naturalness. Atran provides the following example. A no-legged table, suspended from the ceiling is considered a perfectly good table; but a three legged tiger with a prosthetic leg is considered deficient. (Atran, p. 137)

Atran’s model is divided into four taxa in a hierarchy. The taxa in descending order are: folk kingdom, folk life form, folk species, and folk subspecies. This model is diagramed below, with examples, in figure 3.

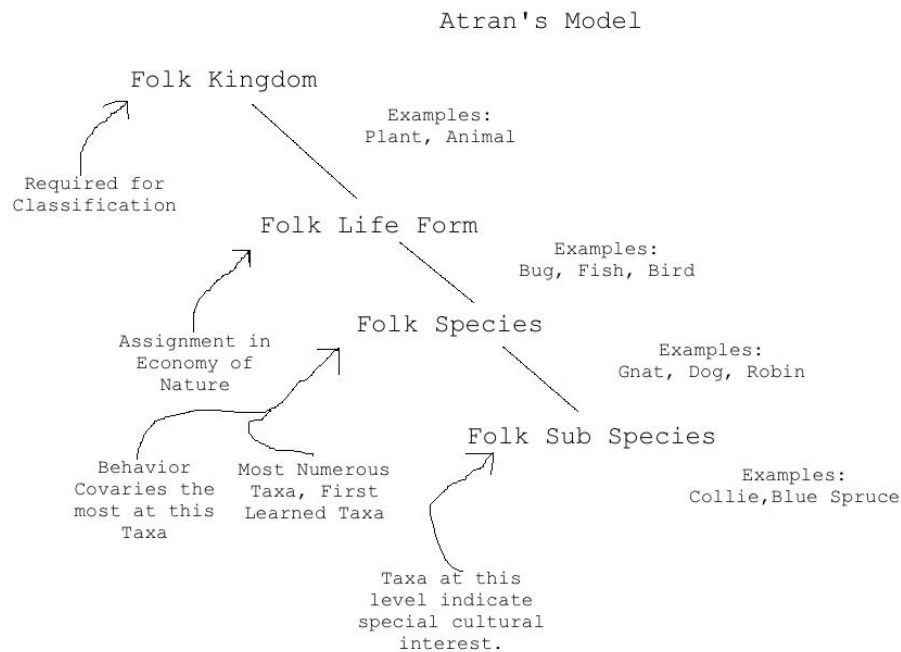


Figure 3: Atran's Model

Atran makes some observations about particular taxa in this system. Of the folk kingdom, he explains, that any observation must be classified into a folk kingdom first if it is to be classified at all. This is a sensible requirement of classification. Additionally, it provides some insight into why scientists are disturbed by the uncertain kingdom of viruses. Of folk life forms, Atran explains

that this class is responsible for the assignment of a classification in the “economy of nature,” (Atran, 1995, p. 142) that is to say, how a particular plant or animal fits into its surroundings. He says of folk species, that they make up the most numerous level in the hierarchy. They are the point at which individual behavior differs the most. Folk species are the first taxa learned by children. They are the most culturally relevant to a people.ⁱⁱⁱ (Atran, 1995, p. 143) This suggests that folk species are akin to the central terms in Hunn’s model, and to the generic level in Berlin’s model. Of folk sub-species, Atran explains, that this is the level of cultural interest. Taxa at this level, for example different varieties of corn, exist because they are of particular interest to a given culture.

While Atran’s model is more similar to Berlin’s than it is to Hunn’s, he shares Hunn’s belief that the field examining ethno-biological classification is too focused on scientific parallels. He states that natural kinds are determined by necessity. (Atran, 1995, p. 164)

Additionally, Atran acknowledges the existence of intermediate taxa that often go unnamed. He provides the example of an intermediate taxon in English with *mouse* and *rat* as children. This taxon accepts no other small rodent. (Atran, 1995, p. 140) Atran believes that, although unnamed, these taxa deserve inclusion in a complete ethno-biological model. This possibility of intermediate taxa is mentioned in Berlin but, because intermediate taxa often go unnamed, Berlin argues against their inclusion as an ethno-biological category. (Berlin, Breedlove, & Raven, 1973, p. 216)

A Primer on Natural Language Processing

Discerning folk taxonomies from a corpus involves sorting through large amounts of data and searching for patterns or similarities in morphology. Procedures making use of large amounts of data are perfectly suited to Natural Language Processing, or NLP for short. Natural Language Processing, also sometimes referred to as Computational Linguistics, is the subfield of Computer Science most concerned with language and its use. There are many tools available to NLP, but the two that I will examine here are: *String Edit Distance* and *WordNet*.

String Edit Distance: finding string similarity

In Computer Science, any arbitrary arrangement of characters is known as a "string." String Edit Distance is a measure of similarity between two strings. The smaller the string edit distance, the more similar the strings. If the string edit distance between two strings is zero, then the two strings in question are identical.

One particular implementation of String Edit Distance is known as "Levenshtein String Distance." This algorithm steps through each pairing of words and scores that pairing. This score is the minimum number of changes that must be made from one string to get to the other. The algorithm understands three operations at any given character, these are: deletion, insertion, and substitution. If any of these three operations is necessary, a point is added to the string edit distance between the two strings. Levenshtein String Distance keeps track of the edit distance of each substring of length n in word a to the corresponding substring of length n in word b . The algorithm can then add this distance to the distance

gained by adding the $n+1$ letter to strings a and b . An example of the computation of Levenshtein String Distance is shown below in figure 4, where the Matukar words for wave, “*lalar*,” and firefly, “*altot*” are compared. The distance between each substring of these two words is shown in their respective cells. For instance might see that the transformation between the substrings ‘*ALT*’ and ‘*LAL*’ can be achieved in two edits, one deletion ‘*T*’ and one addition ‘*L*’.

		<i>A</i>	<i>L</i>	<i>T</i>	<i>O</i>	<i>T</i>
	<i>0</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>
<i>L</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>
<i>A</i>	<i>2</i>	<i>1</i>	<i>2</i>	<i>2</i>	<i>3</i>	<i>4</i>
<i>L</i>	<i>3</i>	<i>2</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>
<i>O</i>	<i>4</i>	<i>3</i>	<i>2</i>	<i>2</i>	<i>2</i>	<i>3</i>
<i>R</i>	<i>5</i>	<i>4</i>	<i>3</i>	<i>3</i>	<i>3</i>	<i>3</i>

Figure 4: Levenshtein String Edit Distance Example

The importance of string similarity may be seen in Berlin’s explanation of the morphology of taxa. Berlin shows that taxa are made up either of primary or secondary lexemes. Primary lexemes are further subdivided into analyzable and un-analyzable groups. (e.g. ‘*crabgrass*’ is analyzable while ‘*grass*’ is not) (Berlin, Breedlove, & Raven, 1973, p. 218) The reason both analyzable primary lexemes and the whole group of secondary lexemes may be analyzed is that they contain embedded words. These morphological similarities provide hints at the underlying order of the folk taxonomy. For example, the secondary lexeme ‘*white rose*’ is a combination of the primary lexeme ‘*rose*’ with the color term ‘*white*’. If we wanted to examine the various varieties of roses in English, we could look for every instance of the word ‘*rose*’ in a complete dictionary and the result would be a list containing all

roses (as well as some noise, such as 'arose'.) This would give us a window into the English folk taxonomic specific children of the taxonomic generic rose.

The above is only possible because we know that English forms binomials in which the second word is rose for its specific rose taxon. The question is, how might we find these analyzable taxa without knowing what any of the language specific patterns are to start? This is where string similarity becomes useful. If string edit distance is run on an entire dictionary, and the most similar words are reported, then words such as “Colorado Spruce” and “Blue Spruce” would be relatively similar due to their second words being identical. Thus, string similarity is a useful tool in an automated taxonomic search.

WordNet: A Semantic Hierarchy of English

The second of the NLP tools of which I make use is WordNet. WordNet is a powerful resource created by Princeton’s Computer Science and Linguistics departments. It may be accessed online at <http://wordnet.princeton.edu>. It contains a relatively comprehensive hand annotated semantic hierarchy for English. WordNet is, in essence, an attempt to provide a solid reference to English’s categorization scheme. English words in WordNet are grouped into sets of “cognitive synonyms,” known as synsets. (Miller, 2011) Synsets are linked together by semantic relations. For example, the synset containing “dog” is a child of the synset containing “domestic animal” and also a child of the synset containing “canine.” Children of the synset containing “dog” include but are not limited to, “puppy,” “poodle,” and “corgi.” A node with a selection of its hypernyms and hyponyms is illustrated in figure 5.

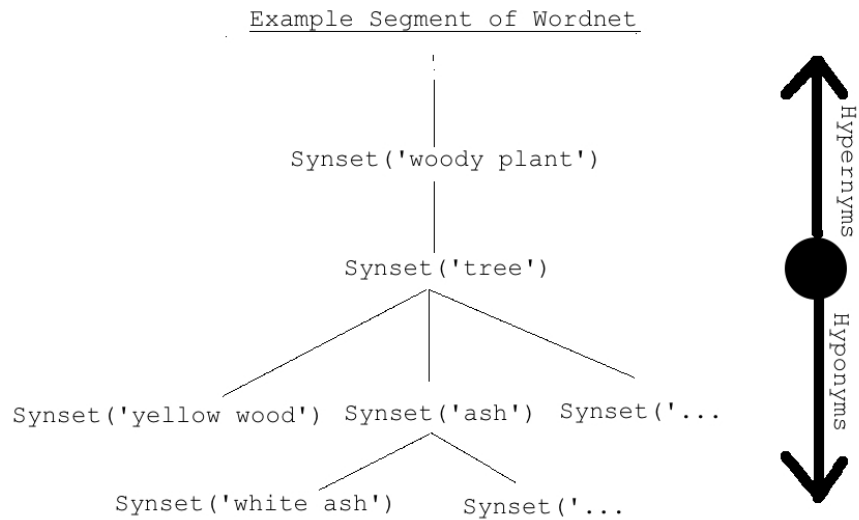


Figure 5: Example Segment of Wordnet

The structure of WordNet closely resembles the hierarchies described by Atran and Berlin. This suggested to me that there might be some way to fit entries in a target language into the English taxonomic tree. Thus, the idea of gloss assisted semantic tagging occurred to me. By using the English gloss for each of a target language's words, I hypothesize that I will be able to tag the words with English semantic fields. I can then walk through the semantic fields and examine the groups for morphological patterns.

One flaw with this approach is that it models the target language onto the English taxonomy, while the points of interest are the target's taxonomy. The hope is that this initial mapping of the target language onto English might provide sets of animals that are similar in English that can later be analyzed for similar morphological qualities in the target language.

It should also be noted that it is only in the best-case scenario that this approach will remove all hand examination of the results. The main intent of this approach is to provide some semantic grouping to an untagged dictionary for the purpose of easing hand analysis afterwards. If the scope of this program is limited to all of the plants and animal words in the dictionary, this should accomplish a categorization of all of the plants and animals in the target language into some more easily understandable format.

Implementation

In this section I will briefly describe the materials and methods I used to leverage the above tools on the Matukar talking dictionary. It should be noted that I was given an XML dump of the dictionary as my corpus. Both of these methods were implemented on Mac OSX using Python 2.7.1. The code for both of these implementations will be available online.^{iv}

String Similarity

For string similarity, I initially implemented Levenshtein String Edit Distance, however; a problem quickly appeared. Levenshtein String Edit Distance does not reward similarity, while it does punish differences. For instance, the string edit distance between “white rose” and “yellow rose” is six, while the distance between “white rose” and “white house” is only two. Words that should have been grouped together were farther apart due to differing length, while words of similar length but differing meanings were being grouped together. This was problematic to say

the least. The results returned by the simple Levenshtein String Edit Distance contained far more noise than they did signal. There are modifications that can be made for Levenshtein String Edit Distance so that it places a reward on similarities between words however; this is more difficult to implement, and it would not necessarily remove the aforementioned problems. It was for this reason that I opted to look into other string similarity algorithms. I found a function, to this purpose, in Python's `difflib` library. This function is `difflib.get_close_matches()`. The help file for `difflib` states that this function implements an advanced version of an algorithm called the "gestalt algorithm," by Ratcliffe and Obershelp^v, to produce similar strings that "look right to humans." This function works by finding the longest subsequence in common between two strings. It then runs the algorithm again on the sequences to the left and right of the previously matched sequences. This alternative sounded promising, and when it was integrated into the program it performed better than the basic string edit distance had, matching fewer sets of words erroneously. This program parsed the Matukar XML file into a dictionary of words, which was then analyzed. Groupings of similar words were generated for each noun in the dictionary. The runtime of this algorithm is relatively fast, taking on the order of a minute or two for the 3045 words in the dictionary.

Gloss Assisted Semantic Tagging

I implemented this method with the use of NLTK, the natural language toolkit, for Python. This method begins by finding all nouns in the dictionary which contain a word in their gloss that is part of a synset a . This synset a is, itself, a child of the synset containing "organism." The intention of this initial step was to collect

all words to which Atran would refer as “living-kinds.” The program then recourses down the hierarchy of synsets, starting from Organism, creating lists of organisms that descend from the current synset. The program stops examining a branch when the current synset has no hyponyms. At each level, when the list of descendant organisms is compiled, if the list is non-empty, it is written to a file so that incremental results may be examined. This was implemented with an object oriented approach with a *NetWalker()* class handling the recursive process and a *NetOrganizer()* class handling the problem of listing descendants. The runtime of this program is rather long as there are many comparisons being made. For the Matukar dictionary it takes about an hour and a half to tag every word in the dictionary for every synset in WordNet that is a child of Organism.

Results of Natural Language Processing

I will discuss and analyze the outputs of these programs, and assess the usefulness of these methods in this section. Both methods had quirks; however, they both demonstrated potential for broader use for future automated analysis of target languages. Sample output for each of these methods may be found in the appendix. Additionally, full output of these programs will be available online.^{vi}

String Similarity: Overview

The use of string similarity as a method of detecting similar lexeme patterns, which should subsequently detect taxonomic groupings, returned some interesting results. There were some instances of success. To begin with, there appeared to be a

fully formed taxonomic set of three birds. The set of similar words to the word for 'chicken' was as follows:

kukurek -chicken
kukurekparpar -hawk (chicken + sound of hawk?)
nubanen kukurek - goose (water + chicken)
kukurek katalun -chicken egg. (chicken + egg)

Figure 6: String similarity of 'kukurek'

With the exception of *kukurek katalun*, 'chicken egg', each of these words represents a different type of bird. This is an exciting result providing evidence in favor of this method. By string similarity alone these words were separated from the entire dictionary. Unfortunately this is the only obvious example of primary/secondary lexeme interaction between classifications of species that I found. This does not mean that the method is unable to pick up on them; it just appears that they may not be present in the Matukar vocabulary, or (even more probably) in the dictionary.

Additional evidence for the utility of this method may be found in the plethora of terms associated with both coconuts and betel plants. In each case the basic words for 'coconut' or 'betel', *niu* and *mariu*, are appended with some other descriptor. (e.g.: *niu patawan*, meaning 'coconut milk'.) For each plant, these terms were grouped into that plant's similar strings. The relevant string edit distance groupings for these words are shown below in Figure 7:

niu ririn - fresh coconut meat remaining in coconut shell after scraping
niu dabin - coconut roots
niu patawan - coconut milk
niu raun - coconut leaf
=====
mariu bag - betel bunch
mariu - betel nut
mariu luwan - betel trunk
mariu digot - betel leaf attachment to tree
mariu sadaro - betel branch (broom)
mariu rau.un - betel leaf

Figure 7: String similarity grouping: niu- 'coconut' and mariu- 'betel'

These groupings do not represent individual species, however, and I have opted not to include them in my analysis of the folk taxonomy of Matukar.

There is also some evidence that, by using this method, the origins of analyzable primary lexemes, in a target language, may be more easily derived. For instance, one Matukar word for 'frog' is *sidar*. The string similarity program returned that this was similar to the Matukar words for both 'blood', *dar*, and 'reef,' *sar*. It is possible that these words are conjoined in some way to create the primary lexeme *sidar*.

Overall there were some promising results for this method; however, due to the relative lack of biological terms in the dictionary, it is difficult to ascertain how effective it is. If there were more diversity in the species elicited for the dictionary, then it would be easier to gauge the effectiveness of this method.

String Similarity: Room for Improvement

While the method of string similarity I used unearthed some interesting patterns, there was still much room for improvement. Some of the below issues are

inherent to this tool, while others have the potential to be mitigated with more advanced techniques. To begin with, this method sorts groups of short words together. These short words, even when very similar in form, often seem to have little to do with each other. Such a grouping may be seen below in Figure 8:

yad - part of a canoe
ya - hole
yau - fire [*paia*]
yan - yellow
dad - buy
bad - pot

Figure 8: Improper Grouping of Short Strings

Aside from a potential relationship between *yau* - 'fire' and *yan* - 'yellow' the other words in this grouping seem unrelated. This occurs because the shortest words have the least opportunity to develop string edit distance. Two three letter words can only be, at most, three string edits apart. This leads to misleading conclusions such as the strings 'cat' and 'sum,' with string edit distances of three, being more similar than the strings 'friend' and 'friendship,' with a string edit distance of four. The former are unrelated, while the latter have the same root. Potential solutions to this problem involve providing rewards to strings with longer similar substrings. For instance, if we decremented the string edit distance for common sub strings then the distance between 'friend' and 'friendship' would be negative two. Such a distance would provide strong evidence for the relatedness of two strings.

A second weakness in string similarity may be seen in the case of binomials with shared descriptors. These descriptors are usually common words. In the output of my program there are many groupings that appear similar to the following in Figure 9:

te dabok - big bilum
nina dabok - big knife
maror dabok - big chief
tamat dabok - big man

Figure 9: Improper Grouping by Binomial Descriptor

These strings were marked as similar due to their shared descriptor *dabok* - 'big.' This would be akin to grouping 'red rose,' 'red fox,' and 'red panda' in English. While these patterns might be interesting, they are outside of our desired results. These errors are an unavoidable byproduct of this method; however, they are usually put into their own groupings and do not impede hand analysis.

Gloss Assisted Semantic Tagging: Overview

The use of WordNet to analyze the glosses of the Matukar dictionary returned interesting results, both promising and problematic. It successfully placed many of the Matukar dictionary entries in their corresponding locations in the English semantic web. This was most often true in the case of plants and animals. I have included the output for the synset 'ant' below in Figure 10:

Synset('ant.n.01')
ror: type of ant (black)
dəm: type of ant (very small, eats sugar)
bakbak: type of ant (black and brown, really big ant...)
kakad: type of ant (big, red ant that goes up tree)
maniŋkal: type of ant (brown, middle sized)
wes: type of ant (black, little ant who bites)

Figure 10: Example Output of NetWalker

The above shows all of the dictionary entries tagged by NetWalker as ants. All of the above entries were tagged correctly. The trigger for categorization into this synset and the synset in question were the same; both were 'ant.' This is not always the case. For instance, in the synset 'insect' we may see, amongst others, the Matukar

words *gab rairai*, ('type of fly,') and *muimui*, ('louse larva.'). These terms were both categorized into the synset 'insects' because NetCrawler identified a string in their gloss, 'fly' and 'larva' respectively, that was an inherited hyponym of insect. Most of the time this method of tagging was sufficient; however, it was not without its flaws.

One problem that appeared in my experiments with this method was that WordNet seems to have included the synset containing 'person' as a hyponym of 'Organism.' While people are certainly organisms, the hyponyms of 'person' in WordNet are societal rolls. This is problematic because the program attempted to tag all of the nouns in the Matukar dictionary with person descriptors such as 'painter' or 'law man.' Even these unintended 'person' related tags were accomplished to some degree of success. For instance the Matukar words for both 'virgin male' and 'virgin female' were tagged under the synset "innocent." This example represents the exception. The noise to signal ratio would have been greatly reduced had 'person' not been included as a hyponym of organism.

Additionally, while browsing the output, I noticed that the Matukar word for 'tilapia' had not found its way into the results tagged with 'fish.' This turned out to be because WordNet categorizes some specific names of animals under the synset "taxa" and not under "organism." I ran the program again, this time with "taxa" as the root, and it returned only one entry all the way to the bottom branch. This was 'tilapia.' I am uncertain whether WordNet has any more words like this, but I am certain that beyond 'tilapia' our analysis of the Matukar dictionary was unaffected.

The most common error, and the only unassailable flaw of this method is improper categorization due to English semantic ambiguity. An example of this is

the improper categorization into the English synset 'gum tree' of the Matukar word *gahu*, which may be translated as 'my gums.' This is a relatively common error and suggests that the output of this method is most useful when checked by hand afterwards.

In the case of Matukar, the output of this program provided all of the same insights as did the string similarity program and more. One piece of information this method detected that the string similarity algorithm missed was the taxonomic class of *is*, the Matukar word for mosquito. When I examined the synset for "mosquito" I noticed that this program had tagged *is*, *is kaduman*, and *is wawak* all as members of this synset. This, in addition to the earlier group (*kukurek*) appears to be a second taxonomic group in Matukar. The reason that string similarity had missed this group was that the element that they all shared in common, *is*, is only two characters long. String similarity did not give appropriate weight to the similar qualities between these terms, as their shared quality was short, and thus passed over them.

I believe that gloss assisted semantic tagging provides an interesting automated means of semantic tagging for any Target English glossed dictionary and seems to produce an understandable hierarchy of organisms in that target language. This could be an invaluable tool to any ethno-biologist. It has a few kinks; however, many of these would be fixable with time, and all of them are recognizable on sight.

Analysis of Output

Matukar seems to have a structured Folk Taxonomy; however, from the data provided in the online talking dictionary I can only find two cases of direct taxa hierarchy. Aside from “*kukurek*,” “*is*,” and their respective descendants the vast majority of the language appears to be at the level of Berlin’s generic taxa. In Hunn’s model, this would suggest that all of the words, save the descendants of the two taxa above, would be central taxa. At first this seems extremely unlikely; however, the purpose of the Matukar dictionary was not to elicit an exhaustive catalogue of their biological terms. Its purpose was to create an initial repository for the language in general. This suggests that the vast majority of living-kind terms elicited were those that were most important to the Matukar people. These relevant terms would be the generic, or central, taxa. A piece of evidence in favor of this explanation is that the vast majority of organism terms found in the dictionary are focused on coconuts, and swine.^{vii} These are both staples of the Matukar way of life and thus would be likely to generate several generics.

One glaring oddity is the absence of life form words, which are hypothesized in both Berlin’s and Atran’s models. Examples of life form words are 'bird,' 'fish,' 'insect,' 'flower.' The only example of a life form word that I was able to find in the dictionary was found in the definition of “bark.” This was “*ai suluṅan*” which literally means 'tree skin.' This suggests that the Matukar for tree is “*ai*”; however, this term was not given its own entry in the dictionary.^{viii}

The results seen here suggest that, in an effort to uncover the folksonomy of Matukar, additional research into the ethno-biology of the Matukar people would be

fruitful. From this initial elicitation few ethno-biological levels are discernable. It would be difficult to continue examination of the Matukar Folk Taxonomy without the ability to elicit additional biological terms, and investigate whether the Matukar people have sets of life forms.

Concluding Thoughts

Future Work

This project has many potential extensions. The string similarity method that I used was more sophisticated than simple Levenshtine String Edit Distance; however, results could be improved further with the utilization of an even more sophisticated string similarity algorithm.

Additional Natural Language Processing tools could be mobilized for this problem. Morphological splitting is a method that, given a training set and a large quantity of words in a target language, attempts to split words into their morphological parts. Morphological splitting seems similar to the way that I use string similarity. Morphological splitting; however, is tuned to search for small strings at the extremes of words. This method could potentially have detected the taxon *is* - 'mosquito' on which my string similarity failed.

The tools that I used can be utilized with the assumption that no large body of literary works exists for the target language. If the researcher had available a large corpus of natural text/speech in the target language, then additional tools

would become available. One example of such a tool is traditional semantic tagging, which attempts to learn the use case for a word by examining copious data.

Bioinformatics tools often provide a suite of web interfaces, and useful visualization tools to researchers. I feel that the methods used in my work with Matukar would scale well to web applications similar to bioinformatics tools such as Basic Linear Alignment Search Tool (BLAST) or ClustalW. These tools could be useful to field researchers who would like some basic automated analysis of a target language.

Conclusion

The study of humanity's categorization of its surrounding is fascinating. That we naturally store our experiences using models for easy recollection is a testament to the efficiency of the human mind. Progress in studies of this area can be easily augmented with several Natural Language Processing techniques. The two techniques examined in this discussion were helpful in making sense of the Matukar Folk Taxonomy and pointing the way for further study.

Bibliography

Atran, S. (1995). Classifying Nature Across Cultures. (E. E. Smith, & D. N. Osherson, Eds.) *An Invitation to Cognitive Science*, III.

Berlin, B., Breedlove, D. E., & Raven, P. H. (1973). General Principles of Classification and Nomenclature in Folk Biology. *American Anthropologist*, 75, 214-242.

Fellbaum, C. (1998). *Wordnet: An Electronic Lexical Database*. (C. Fellbaum, Ed.) Bradford Books.

Gluck, J. (2011). *NLP Assisted Analysis of Folk Taxonomy*. Swarthmore: Self.

Harrison, K. D., Anderson, G. D., & Mathieu-Reeves, D. (2010, 1 1). *About The Dictionary*. Retrieved 5 4, 2011, from Matukar Online Talking Dictionary: <http://matukar.swarthmore.edu/about.php>

Hunn, E. (1982). The Utilitarian Factor in Folk Biological Classification. *American Anthropologist*, 84, 830-847.

Miller, G. A. (2011). *Princeton University*. Retrieved 5 4, 2011, from Word Net: <http://wordnet.princeton.edu/>

The Central Intelligence Agency. (2009). *CIA-The World Fact Book*. Retrieved 5 4, 2011, from Central Intelligence Agency:

<https://www.cia.gov/library/publications/the-world-factbook/geos/pp.html>

Appendix

Map of Matukar



(Harrison, Anderson, & Mathieu-Reeves, 2010)

Example Output of String Similarity

tim: air
tim: wind
tidom: night
ti: no

nub yahai: waterfall
numau tahaik: five
nub narman: Water from yesterday
i yakai: he goes (but...)
ab yabi: S/he makes a house
nub wananan: hot water
nub koraman: puddle

kukurek: chicken
kukurekparpar: hawk
nubanen kukurek: goose
kukurek katalun: chicken egg

se paiin: paternal grandmother
sise paiin: old woman
sileŋ paiin: laughing woman
paiin: woman
kol paiin: female cousin
ham paiin: your wife
bagebage paiin: grandmother
ŋahau paiin: my wife
i wau paiin: my daughter-in-law
i wam paiin: your daughter-in-law

raurau uyan: Hello
garmaurau.un: my hair
abaŋ uyan: good day
garmauraun: my hair
mariu luwan: betel trunk
nal uyan: good day
fud uyan: good banana

Example Successful Output of NetWalker

```
=====
Synset('arthropod.n.01')
ror: type of ant (black)
kasaromrom: type of spider (lives in house)
dəm: type of ant (very small, eats sugar)
ləd: louse egg
is kaduman: mosquito larva
katabebe: spider
is wawak: mosquito (big)
bakbak: type of ant (black and brown, really big ant, goes up tree)
kabob: butterfly
altot: firefly
kalambu: mosquito net
kaiya: termites
alili: centipede
kakad: type of ant (big, red ant that goes up tree)
maniḡkal: type of ant (brown, middle sized)
is: mosquito
teratettet: type of insect
wes: type of ant (black, little ant who bites)
ut: louse
degadəg: cockroach
gab rairai: type of fly (big, blue)
muimui: louse larva
bukabuk: mosquito bite

=====
Synset('arachnid.n.01')
kasaromrom: type of spider (lives in house)
katabebe: spider

=====
Synset('spider.n.01')
kasaromrom: type of spider (lives in house)
katabebe: spider

=====
Synset('centipede.n.01')
alili: centipede
```

Example Improper Output of NetWalker

=====
Synset('producer.n.02')
mariu pidin: wood from betel nut tree
mariu digot: betel leaf attachment to tree
uləp: rope circle used for climbing trees (goes around feet)
ai suluŋan: bark (lit. tree skin)
nyat: hook for getting something from trees
tabe: brain, noodles, something inside of a rotten tree
pat: stone [(si)ton]

=====
Synset('film_maker.n.01')
pat: stone [(si)ton]

=====
Synset('architect.n.01')
kabakabman: eye white (possessed)
pat: stone [(si)ton]
kabakab: white

=====
Synset('maker.n.01')
laŋalaŋ tatuan: railing post
bag: post

=====
Synset('manufacturer.n.02')
laŋalaŋ tatuan: railing post
bag: post

i Our modern taxonomies may be non-biological in nature, because our surroundings no longer call for biological categorization. One example of a non-biological folk category would be the "chick flick."

ii Map of area attached in appendix

iii Atran shows this with an explanation of how children in the western world recall folk species the most quickly only in cases of mammals. When a non-mammal was elicited, the children produced folk life form terms.

iv The code will be made available at

<http://www.sccs.swarthmore.edu/users/12/jgluck/Files/Linguistics/Matukar/src>

It should be noted that the code will not work without NLTK having been installed.

v A detailed explanation of this algorithm may be found at

<http://drdobbs.com/article/print?articleId=184407970&siteSectionName=>

-
- vi *The output of these programs will be made available at <http://www.sccs.swarthmore.edu/users/12/jgluck/Files/Linguistics/Matukar/Results>*
- vii *These terms were not individual specie terms, they were terms for parts of a coconut tree, or for counting swine.*
- viii *I have since hand checked the XML dump of the dictionary, and found that the word “ai” is included with the gloss “wood.” This gloss did not trigger inclusion in Organism’s hyponyms, because wood is an object, not an organism.*