

Practical Python

7: Basic machine Learning with NLTK



Dafydd Gibbon

First South African Workshop in Digital Humanities
North-Western University, Potchefstroom, SA
2015-04-04 to 2015-04-05

ML with NLTK

- ML types:
 - Naive Bayes
 - Maximum Entropy / Logistic Regression
 - Decision Tree
 - SVM (?)
- NLTK-Trainer
 - `https://github.com/japerk/nltk-trainer`
 - command line scripts
 - train custom models
 - analyse corpora
 - analyse models against corpora

Machine Learning – supervised learning

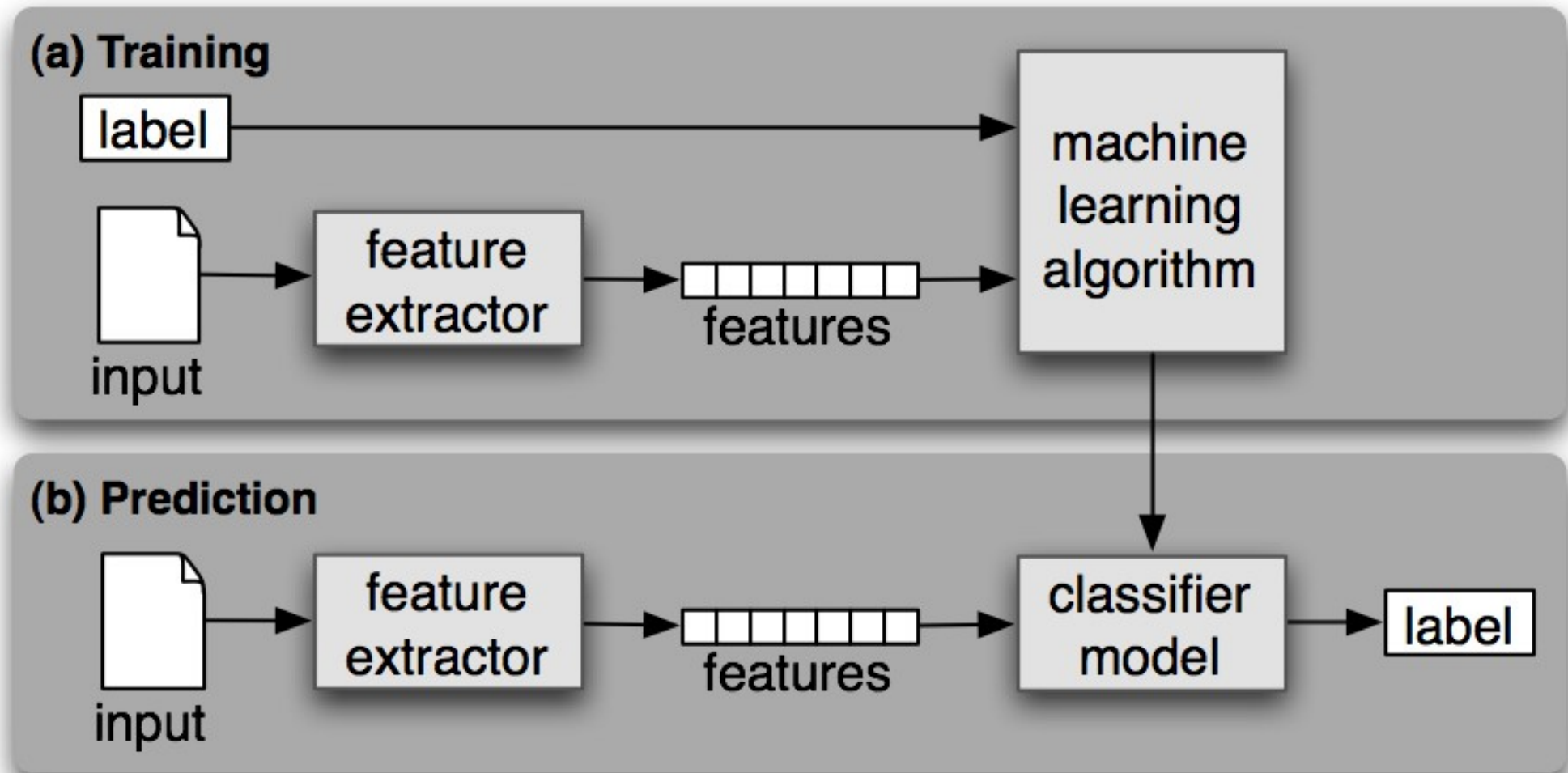


Figure from NLTK book.

A textbot with NLTK (unsupervised learning) (and its limitations):

```
import nltk

def generate_model(cfdist, word, num=15):
    for i in range(num):
        print word,
        word = cfdist[word].max()

# Get data
text = nltk.corpus.genesis.words('english-kjv.txt')
# Define probabilities between neighbouring words:
bigrams = nltk.bigrams(text)
cfd = nltk.ConditionalFreqDist(bi
# Predict
generate_model(cfd, 'living')
```

Suggestion:
Solve the circularity
problem of the textbot.

Simple gender-feature name classifier

```
#!/home/gibbon/miniconda2/bin/python
# Simple gender classifier from NLTK book, Ch. 6 (edited)

import nltk
from nltk.corpus import names
import random

# Define a (kind of) feature extractor
def gender_features(word):
    return {'last_letter': word[-1]}

# Get data
male = [(name, 'male') for name in names.words('male.txt')]
female = [(name, 'female') for name in names.words('female.txt')]
labeled_names = male + female
random.shuffle(labeled_names)

# Get feature sets, partition into training and test sets
featuresets = [(gender_features(n), gender) for (n, gender) in
labeled_names]
train_set = featuresets[500:]
test_set = featuresets[:500]
```

Simple gender-feature name classifier

Training

```
classifier = nltk.NaiveBayesClassifier.train(train_set)
```

Testing

```
accuracy = nltk.classify.accuracy(classifier, test_set)
```

Results

```
print "train:test set sizes:",len(train_set),':',len(test_set)
```

```
print "Accuracy with", len(test_set), "items:", accuracy
```

```
print "Most Informative Features: (likelihood ratios per  
feature):"
```

```
classifier.show_most_informative_features(5)
```

Predicting for single items

```
testlist = ['Robert', 'Anna', 'Neo', 'Trinity']
```

```
maxlen = max([len(x) for x in testlist])
```

```
for x in testlist:
```

```
    gender = classifier.classify(gender_features(x))
```

```
    print ('"' + x + '":').ljust(maxlen+4), gender
```

Learning distances between languages

- The DistGraph web site describes and demonstrates the discovery of similarities and differences (interpreted as ‘distances’) between languages of the Ivory Coast:
 - <http://wwwhomes.uni-bielefeld.de/gibbon/DistGraph/>
- The data are consonant systems of Ivory Coast languages, harvested from atlases of Ivory Coast languages dating back to fieldwork done in the late 1970s and early 1980s.
- Consonant systems are chosen
 - because of their relative ease of identification in comparison to other phonological properties such as vowels or tones,
 - because analysis of lexical similarities (characteristic of many similar studies) presupposes phonological analysis as the main basis of lexical similarities.

A complex Python web application:
Learning distances between languages

Learning distances between languages: workflow

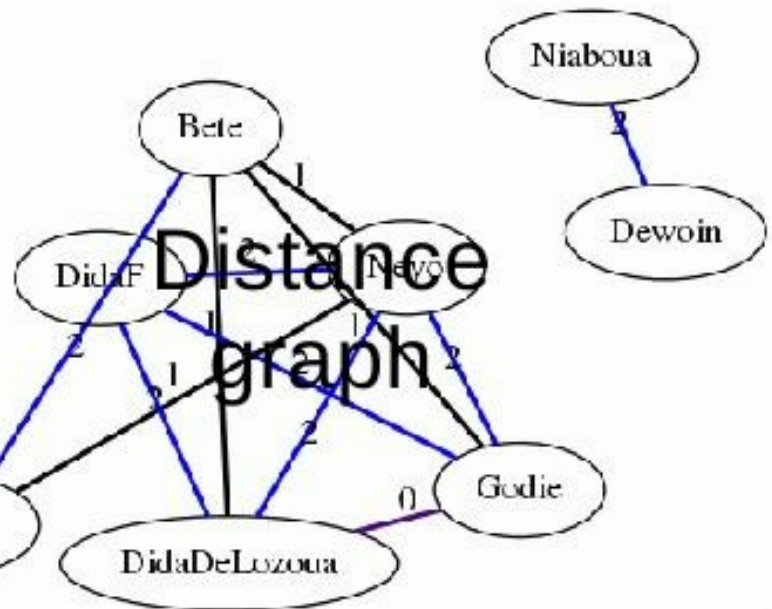


Preprocessing:

1. scan/type
2. correct
3. create CSV

Bete	p	t	c	k	kp	kw	b	d	C	g	sb	f	s	v	z	B	i	j	x	w	m	n	j	N	Nw																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																									</
------	---	---	---	---	----	----	---	---	---	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---	----	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	----

Levenshtein Edit Distance comparison



GraphViz

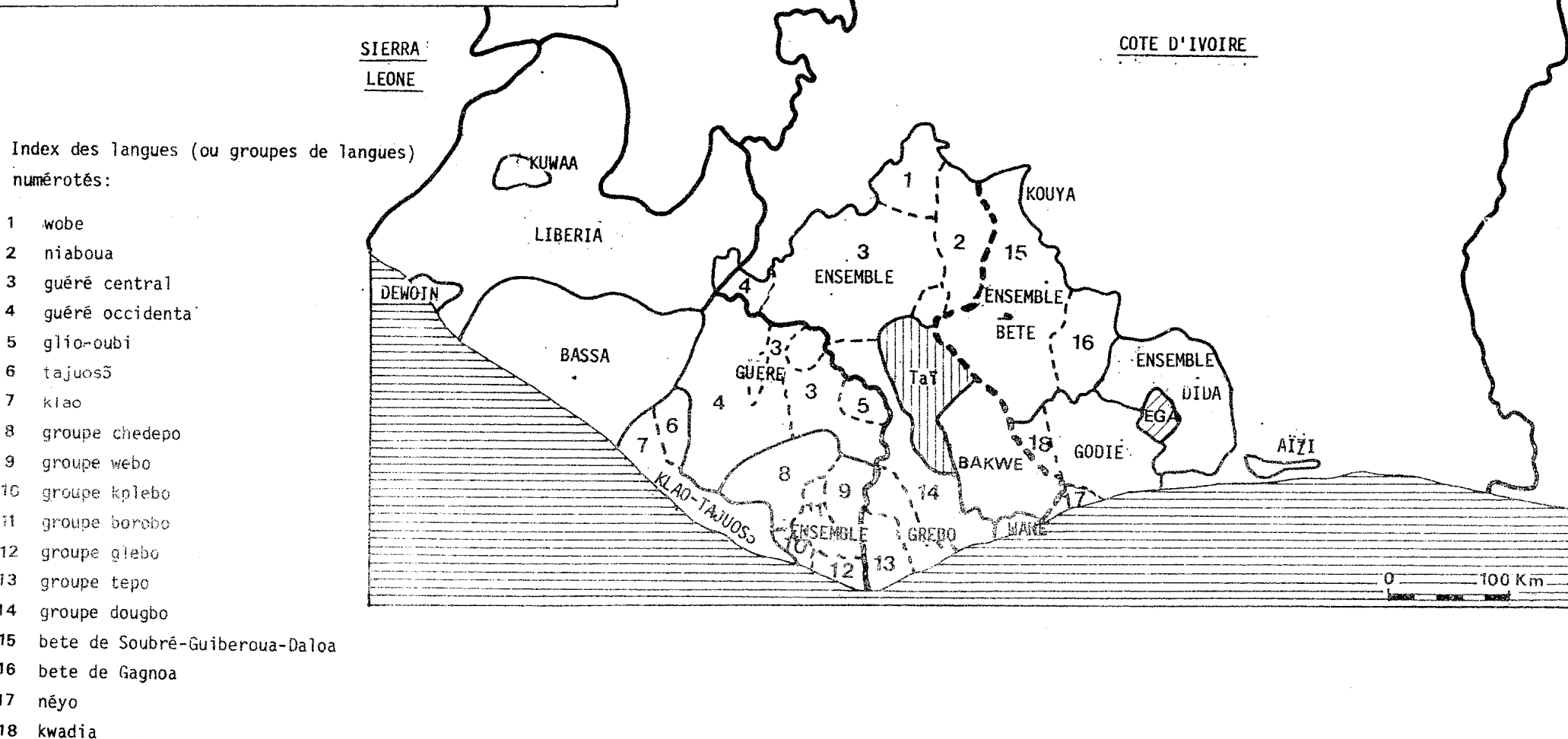
Distance matrix

Bete	0	1	2	1	1	3	10	6	9	11	8	4	4	7	11	8	12	9	6
Godie	1	0	3	2	0	2	11	5	10	12	9	3	3	8	12	8	13	10	7
Koyo	2	3	0	1	3	3	12	8	9	11	8	4	4	9	13	8	12	9	6
Nevo	1	2	1	0	2	2	11	7	8	10	7	3	3	8	12	7	11	8	5
DidiDeLuzoua	0	1	2	0	2	1	11	5	10	12	9	3	3	8	12	8	13	10	7
DidiF	1	0	3	2	0	2	11	5	10	12	9	3	3	8	12	8	13	10	7
Wibe	3	11	8	7	1	1	12	8	9	11	8	4	4	9	13	8	12	9	6
Guere	6	9	8	7	8	7	11	11	11	11	8	4	4	9	13	10	11	10	7
Kuho	6	10	8	7	8	7	11	11	11	11	8	4	4	9	13	10	11	10	7
Cesepo	11	12	11	10	12	10	6	11	4	0	3	9	11	10	5	13	8	11	8
Kiao	8	9	8	7	9	7	4	8	3	3	0	8	8	11	8	4	10	7	8
Niaboua	4	3	3	0	2	2	11	5	10	12	9	3	3	8	12	8	13	10	7
Dewoin	4	3	3	0	2	2	11	5	10	12	9	3	3	8	12	8	13	10	7
Bassia	7	8	7	8	7	8	11	11	11	11	8	4	4	9	13	10	11	10	7
Grebo	11	12	11	10	12	10	6	11	4	0	3	9	11	10	5	13	8	11	8
Teso	8	9	8	7	9	7	4	8	3	3	0	8	8	11	7	0	12	7	8
KoumaLiberia	12	13	12	11	13	11	13	14	10	11	13	10	14	12	19	17	12	0	25
SemeHoumVota	9	10	9	8	10	10	11	11	8	8	7	7	9	8	10	7	15	0	5
AwaCdi	6	7	6	5	7	7	12	10	9	11	8	8	8	9	11	8	14	5	0

Learning distances between languages: hand-drawn language map

- Division entre langues orientales et langues occidentales
- Frontière entre ensembles de langues
- - - Frontière entre langues ou groupes de langues
- Frontière entre pays
- ▨ Enclave non-kru
- ▤ Réserve de Taï (zone inhabitée)

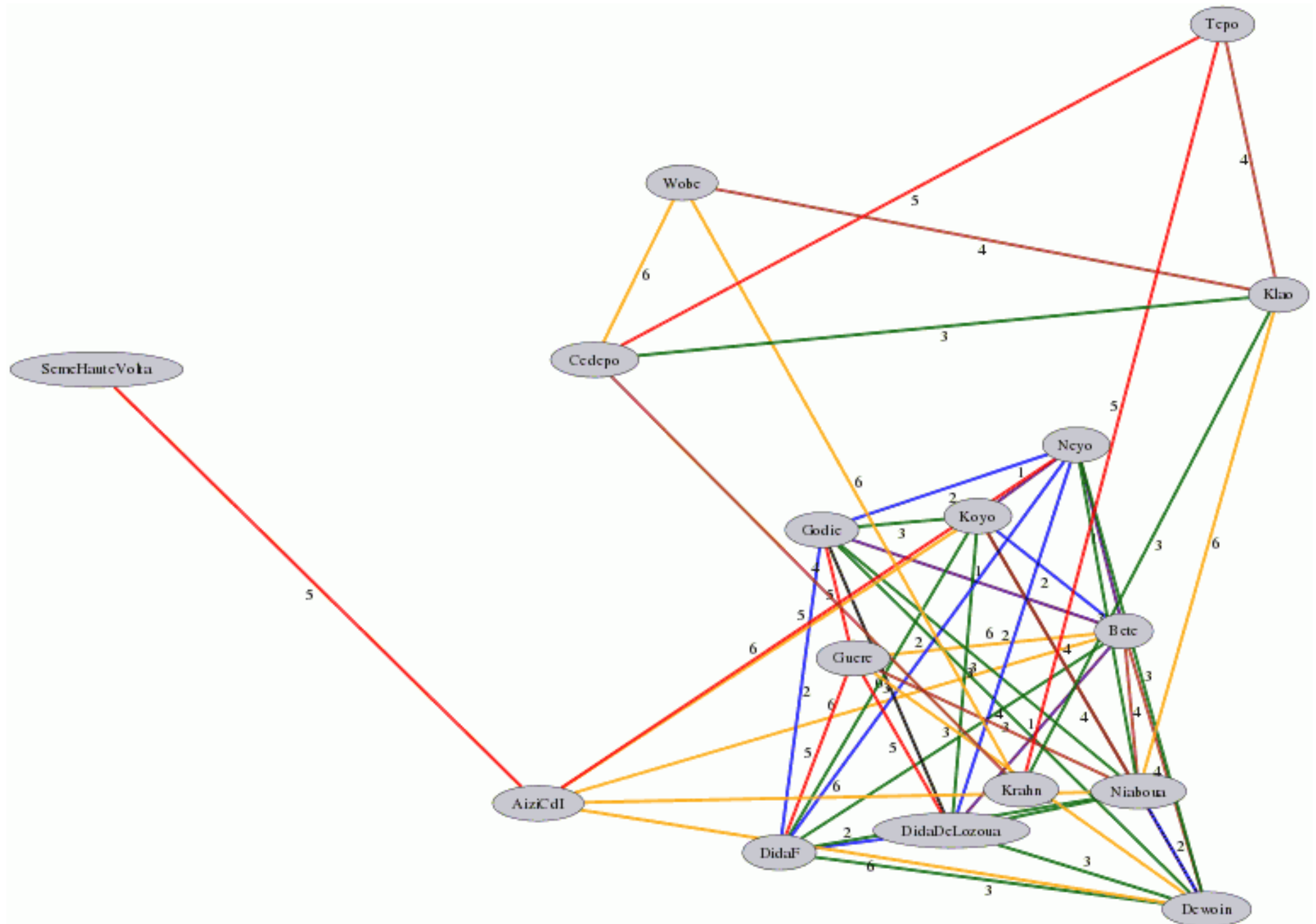
(les groupements de langues pour la partie libérienne de la carte ont été fournis par M. John DUTSMAN.)



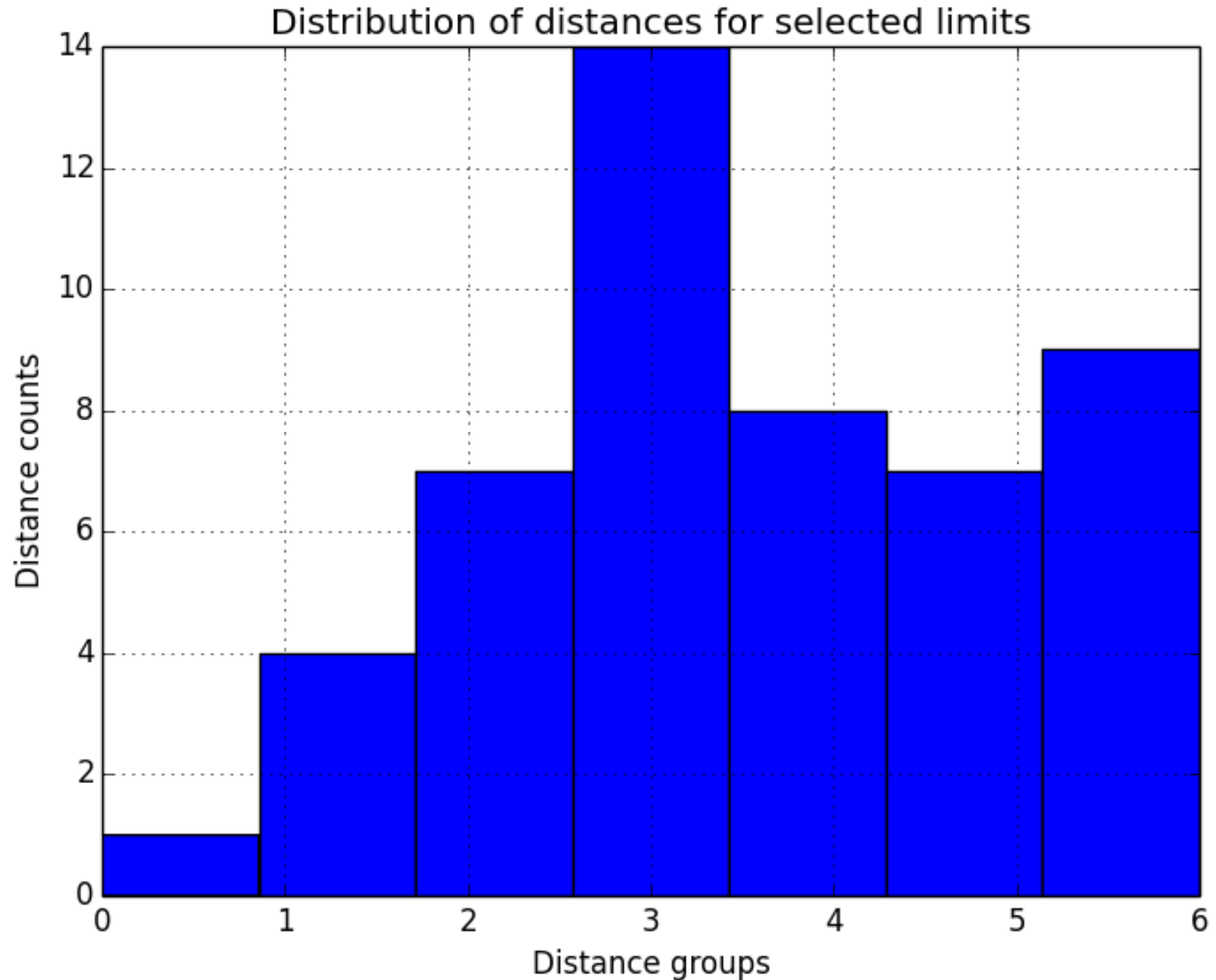
Learning distances between languages: Kru consonant data

Bete;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;_ ;f;s;_ ;v;z;_ ;_ ;_ ;B;_ ;l
Godie;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;_ ;_ ;B;
Koyo;p;t;c;k;kp;kw;kj;b;d;C;_ ;g;gb;_ ;f;s;_ ;v;z;_ ;_ ;_ ;B;_ ;
Neyo;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;_ ;f;s;_ ;v;z;_ ;_ ;_ ;B;_ ;l
DidaDeLozoua;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;
DidaF;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;_ ;_ ;B;
Wobe;p;t;c;k;kp;kw;_ ;b;d;C;_ ;_ ;gb;_ ;f;s;_ ;_ ;_ ;_ ;_ ;_ ;_ ;
Guere;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;_ ;_ ;B;D
Krahn;p;t;c;k;_ ;kw;_ ;b;d;C;_ ;_ ;gb;_ ;f;s;_ ;_ ;_ ;_ ;_ ;_ ;_ ;l
Cedepo;p;t;c;k;kp;kw;_ ;b;d;C;_ ;_ ;gb;_ ;f;s;_ ;_ ;_ ;_ ;h;_ ;_ ;
Klao;p;t;c;k;kp;kw;_ ;b;d;C;_ ;_ ;gb;_ ;f;s;_ ;_ ;_ ;_ ;_ ;_ ;_ ;l
Niaboua;p;t;c;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;_ ;_ ;E
Dewoin;p;t;_ ;k;kp;kw;_ ;b;d;C;_ ;g;gb;gw;f;s;_ ;v;z;_ ;_ ;_ ;B;
Bassa;p;t;c;k;kp;_ ;_ ;b;d;C;dj;g;gb;_ ;f;s;_ ;v;z;_ ;h;hw;B;
Grebo;p;t;c;k;kp;_ ;_ ;b;d;C;_ ;g;gb;_ ;f;s;_ ;_ ;_ ;h;hw;_ ;_ ;
Tepo;p;t;c;k;_ ;kw;_ ;b;d;C;_ ;g;gb;_ ;f;s;_ ;_ ;_ ;h;_ ;_ ;_ ;l
Kuwaaliberia;p;t;_ ;k;kp;kw;_ ;b;d;C;_ ;_ ;_ ;f;s;_ ;_ ;_ ;_ ;
SemeHauteVolta;p;t;c;k;kp;_ ;_ ;b;d;C;_ ;g;gb;_ ;f;s;S;v;_ ;_ ;
AiziCdI;p;t;c;k;kp;_ ;_ ;b;d;C;_ ;g;gb;_ ;f;s;S;v;z;Z;_ ;_ ;_ ;

Learning distances between languages: Hamming Distances



Learning distances between languages: distance histogram

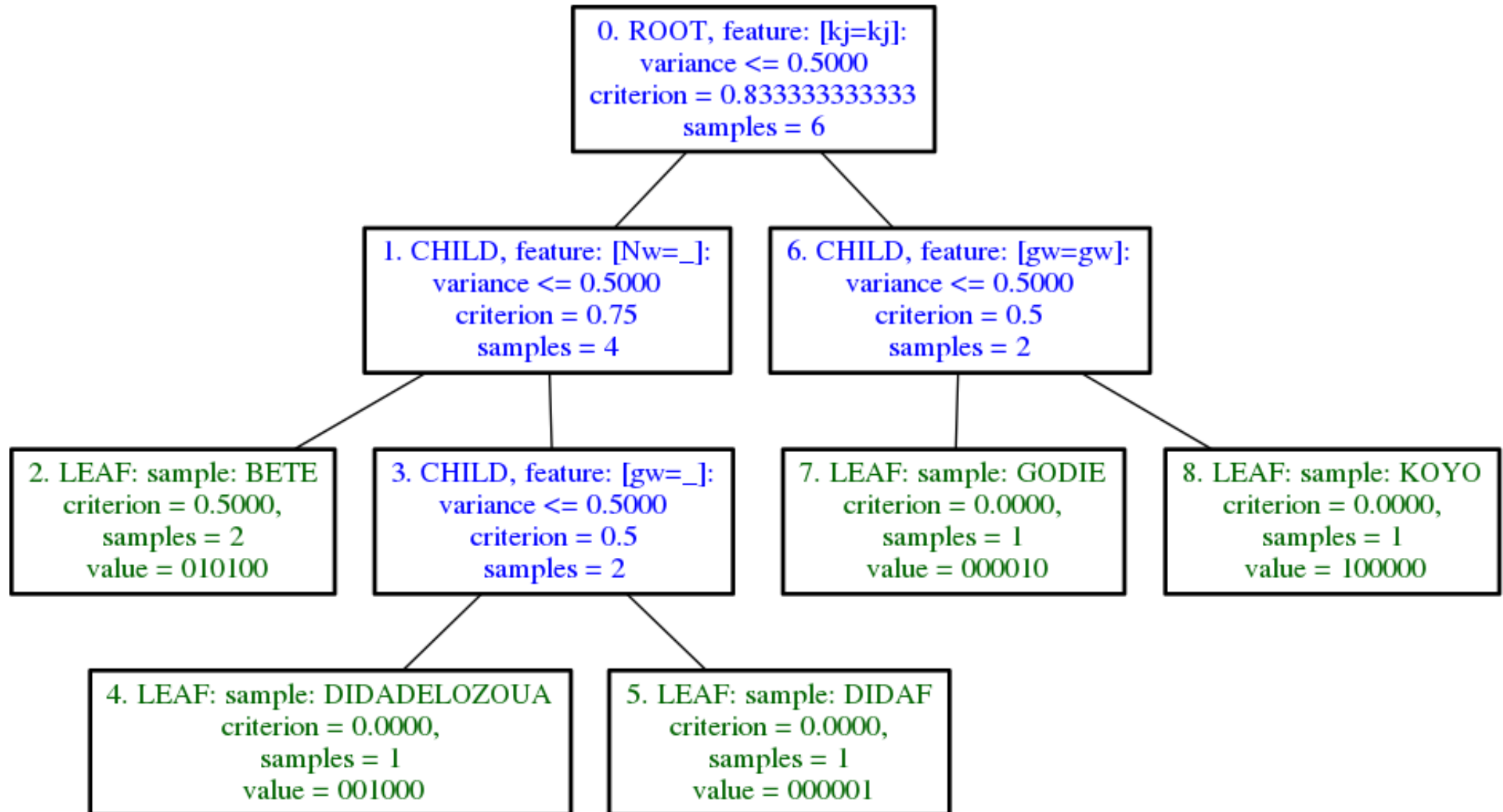


Learning distances between languages: language classification

- Machine Learning:
 - First determine the relative importance of features
 - Then apply a classifier to group the languages in terms of selected features
- Work in progress ... :)

Learning distances between languages: feature importance

Induction of Decision Trees



Learning distances between languages: Python implementation

- Resource: Legacy Ivory Coast Language Atlas
- User interaction via:
 - HTML input form
 - parameter settings
 - CSV table pasting
 - Common Gateway Interface (CGI)
 - Server-side Python application
 - HTML output:
 - Text, Statistics
 - Images linked to graph
- User can download graphs and copy text and statistics
- Processing:
 - CSV transformation to Python list of lists
 - Pairwise application of Levenshtein Edit Distance or Hamming Distance
 - Construction of triangular distance matrix
 - Transformation of matrix into distance triples
 - Conversion of distance triples into GraphViz network code
 - System call to GraphViz dot application for graph generation
 - Graph storage

End of Unit 7