

Probabilistic models in Phonology

John Goldsmith
University of Chicago
Conference on Phonological Variation:
The Case of French
Tromsø: CASTL
August 24, 2005

Web link

For more information on software discussed
here, go to:

<http://humfs1.uchicago.edu/~jagoldsm/Tromso/index.htm>

Probabilistic phonology

Why a phonologist should be interested in probabilistic tools for understanding phonology, and analyzing phonological data...

- Because probabilistic models are very powerful, and can tell us much about data even without recourse to structural assumptions, and
- Probabilistic models can be used to teach us about phonological structure.

The two parts of today's talk will address each of these.

Automatic learning of grammars

Automatic learning of grammars: a conception of what linguistic theory *is*.

Automatic learning techniques:

- In some respects they teach us *more*, and in some respects they teach us *less*, than non-automatic means.
- Today's talk is a guided tour of some applications of known techniques to phonological data.

Probabilistic models

- Are well-understood mathematically;
- Have powerful methods associated with them for learning parameters from data;
- Are the ultimate formal model for understanding *competition*.

Essence of probabilistic models:

- Whenever there is a choice-point in a grammar, we must assign degrees of *expectedness* of each of the different choices.
- And we do this in a way such that these quantities add up to 1.0

Frequencies and probabilities

- **Frequencies** are numbers that we observe (or count);
- **Probabilities** are parameters in a theory.
- We can set our probabilities on the basis of the (observed) frequencies; but we do not *need* to do so.
- We often do so for one good reason:

Maximum likelihood

- A basic principle of empirical success is this:
 - Find the probabilistic model that assigns the highest probability to a (pre-established) set of data (observations).
- Maximize the probability of the data.

Brief digression on Minimum Description Length (MDL) analysis

- Maximizing the probability of the data is not an entirely satisfactory goal: we also need to seek economy of description.
- Otherwise we risk *overfitting* the data.
- We can actually define a better quantity to optimize: this is the *description length*.

Description Length

- The description length of the analysis A of a set of data D is the sum of 2 things:
 - The *length* of the grammar in A (in “bits”);
 - The *(base 2) logarithm of the probability assigned to the data D , by analysis A , times -1* (“log probability of the data”).
- When the probability is high, the “log probability” is small; when the probability is low, the log probability gets large.

MDL (suite)

- If we aim to minimize the sum of the description length (= length of the grammar, as in 1st generation generative grammar) + log probability (data), then we will seek the best overall grammatical account of the data.

Morphology

- Much of my work over the last 8 years has been on applying this framework to the discovery of morphological structure.
- See <http://linguistica.uchicago.edu>
- Today, though: phonology.

Assume structure?

- The standard argument for assuming structure in linguistics is to point out that there are empirical generalizations in the data that cannot be accounted for without assuming the existence of the structure.

- Probabilistic models are capable of modeling a great deal of information without assuming (much) structure, and
- They are also capable of *measuring* exactly how much information they capture, thanks to information theory.
- Data-driven methods might be especially of interest to people studying dialect differences.

Simple segmental representations

- “Unigram” model for French (English, etc.)
- Captures only information about segment frequencies.
- The probability of a word is the product of the probabilities of its segments.
- Better measure: the complexity of a word is its average log probability:

$$\frac{1}{length(W)} \sum_{i=1}^{length(W)} -\log_2 prob(w_i)$$

Let's look at that graphically...

- Because *log probabilities* are much easier to visualize.
- And because the log probability of a whole word is (in this case) just the sum of the log probabilities of the individual phones.

Add (1st order) conditional probabilities

- The probability of a segment is conditioned by the preceding segment.
- Surprisingly, this is mathematically equivalent to *adding* something to the “unigram log probabilities” we just looked at: we add the “mutual information” of each successive phoneme.

$$MI(pq) = \log \frac{prob(pq)}{prob(p)prob(q)}$$

Mutual information

$$MI(pq) = \log \frac{prob(pq)}{prob(p)prog(q)}$$

Weighted mutual information

$$WMI(pq) = prob(pq) \log \frac{prob(pq)}{prob(p) prog(q)}$$

Complexity = average log probability

- Find the model that makes this equation work the best.
- Rank words from a language by complexity:
 - Words at the top are the “best”;
 - Words at the bottom are borrowings, onomatopoeia, etc.

- The pressure for nativization is the pressure to rise in this hierarchy of words.
- We can thus define the direction of the phonological pressure...

Nativization of a word

- Gasoil [gazojl] or [gazɔl]
- Compare average log probability (bigram model)
 - [gazojl] 5.285
 - [gazɔl] 3.979
- This is a huge difference.
- Nativization *decreases the average log probability of a word.*

Phonotactics

- Phonotactics include knowledge of 2nd order conditional probabilities.
- Examples from English...

1 stations
2 hounding
3 wasting
4 dispensing
5 gardens
6 fumbling
7 telescience
8 disapproves
9 tinker
10 observant
11 outfitted
12 diphtheria

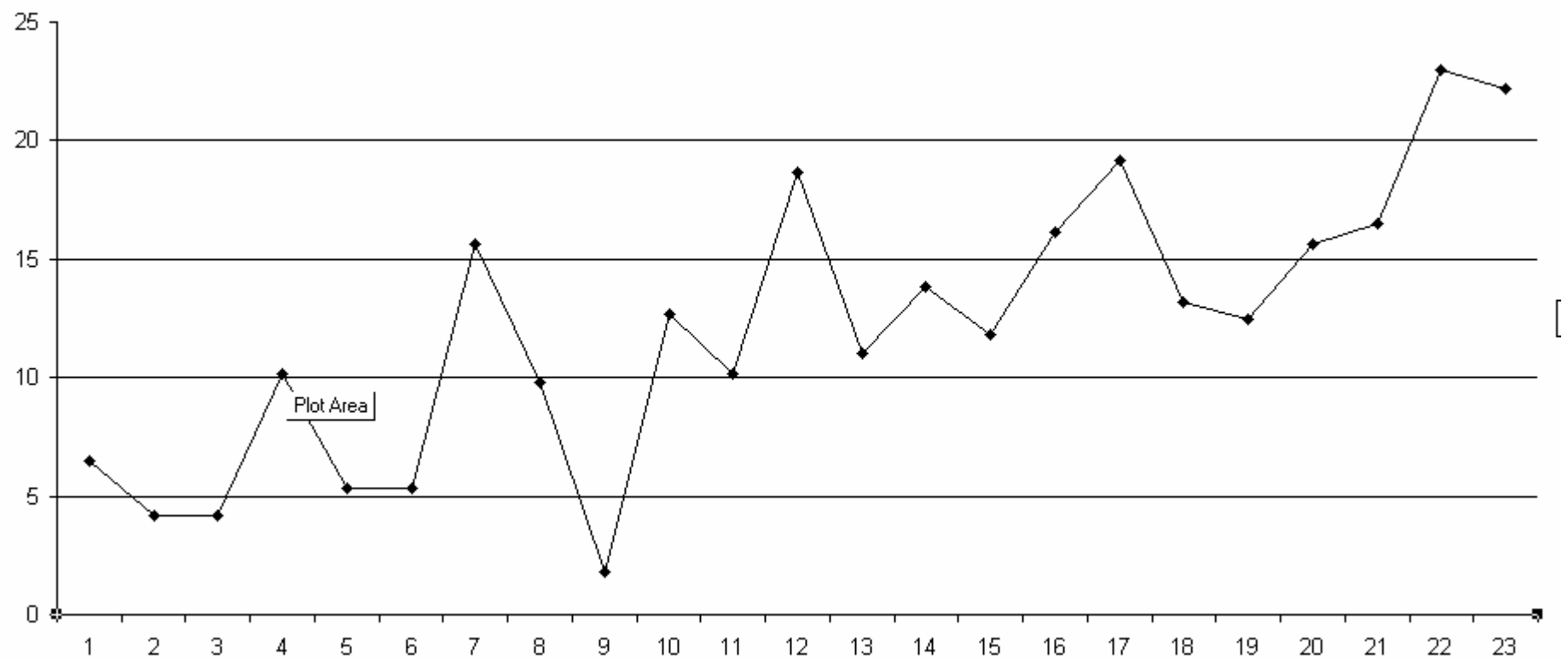
13 voyager
14 schaffer
15 engage
16 Louisa
17 sauté
18 zigzagged
19 Gilmour
20 Aha
21 Ely
22 Zhikov
23 kukje

But speakers didn't always agree. The biggest disagreements were:

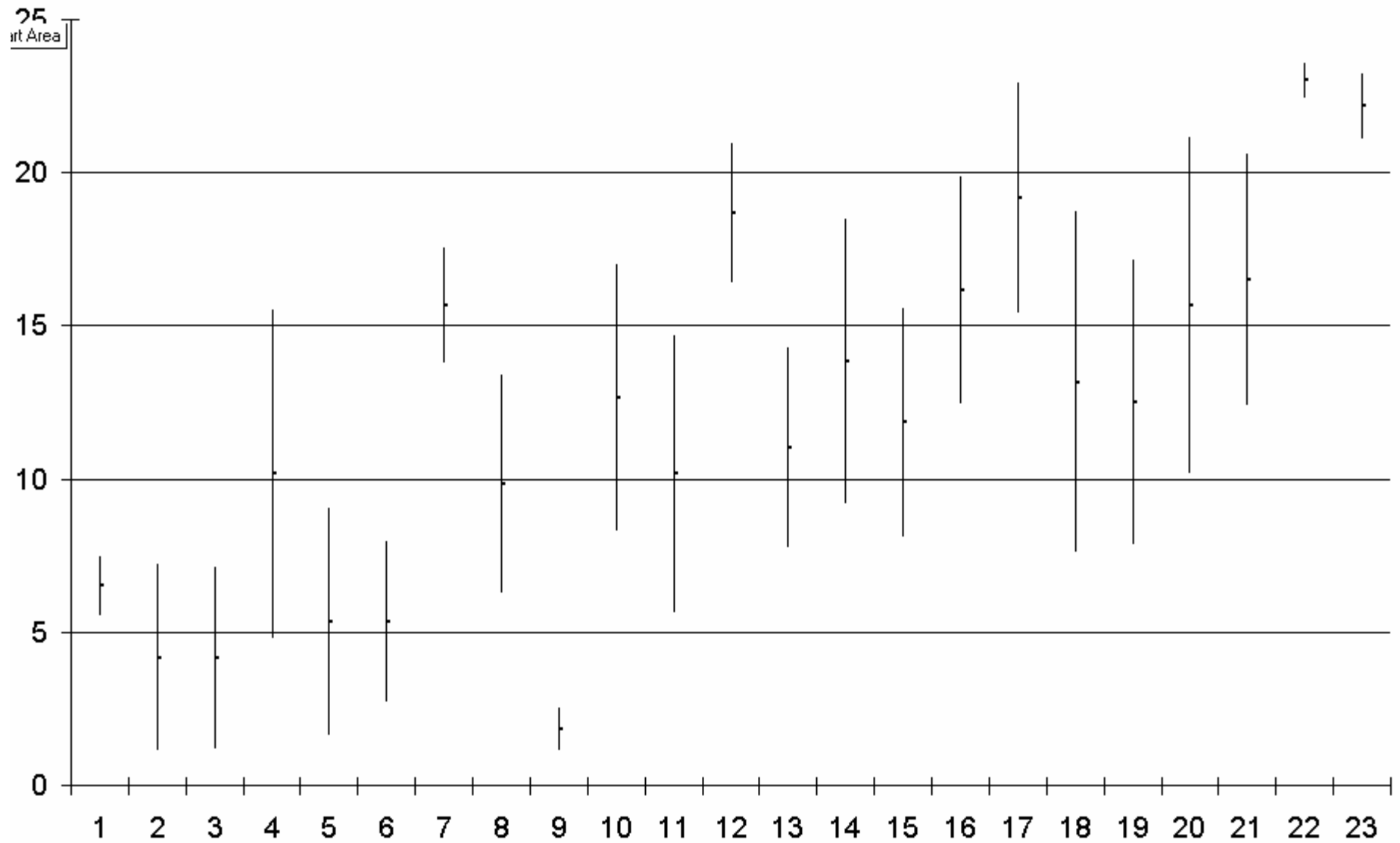
People liked this better than computer:
tinker

Computer liked this better than people:
dispensing, telesciences, diphtheria, sauté

Here is the average ranking assigned by six speakers:



and here is the same score, with an indication of one standard deviation above and below:



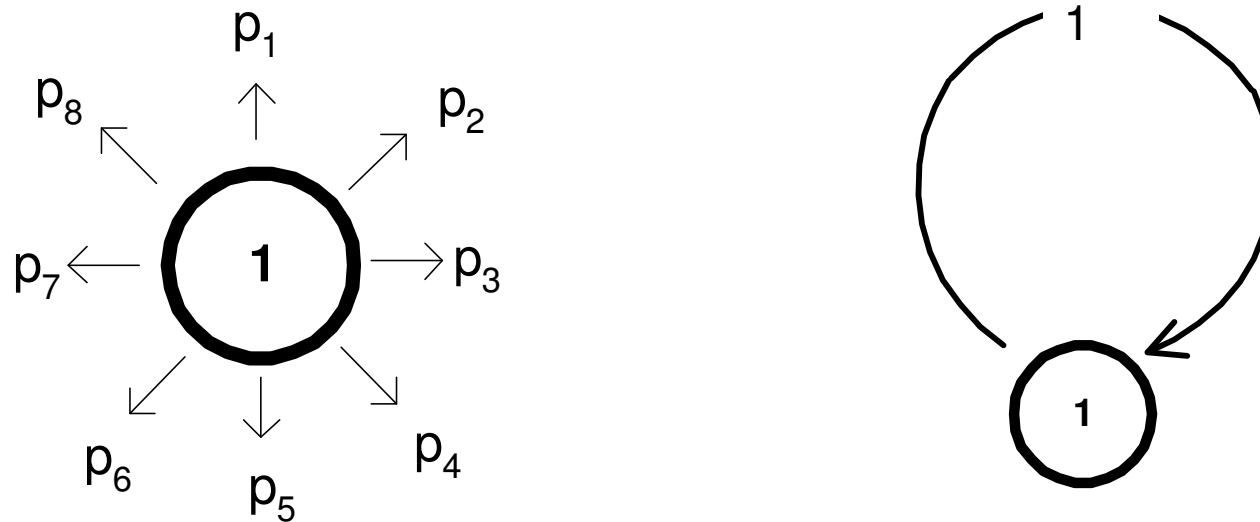
Part 2: Categories

- So far we have made no assumptions about categories.
- Except that there are “phonemes” of some sort in a language, and that they can be counted.
- We have made no assumption about phonemes being sorted into categories.

Emitting a phoneme

- We will look at models that do *two* things at each moment:
- They *move from state to state*, with a probability assigned to that movement; and
- They emit a symbol, with a probability assigned to emitting each symbol.
- The probability of the entire path is obtained by multiplying together all of the state-to-state transition probabilities, *and* all of the emission probabilities.

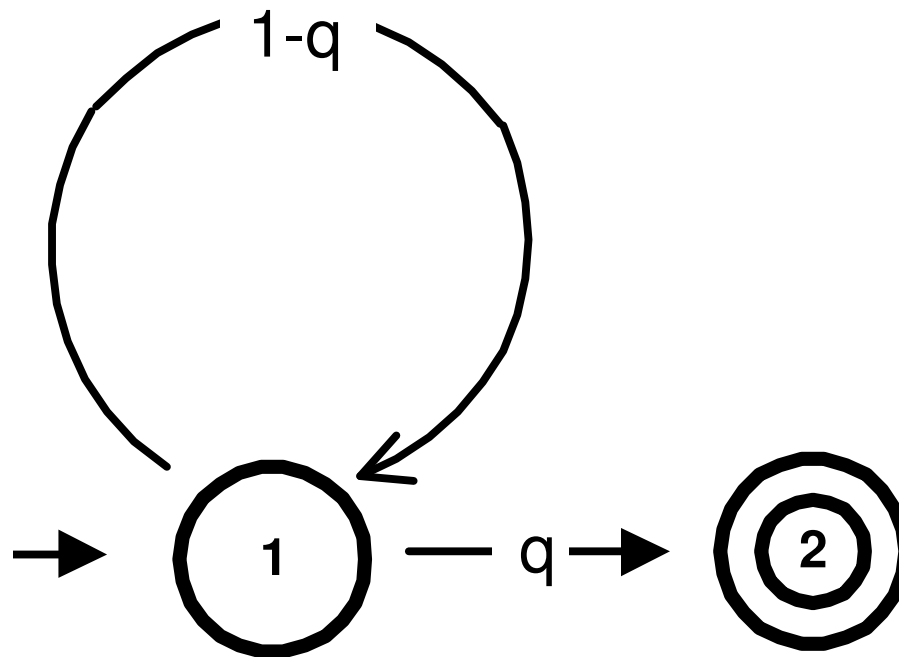
Simplest model for producing the strings of phonemes observed for a corpus (language)



To emit a sequence p_1p_2 and stop, there is only one way to do it:
Pass through state 1 *twice*, then stop.
The steps will “cost”: $p_1 * p_2$

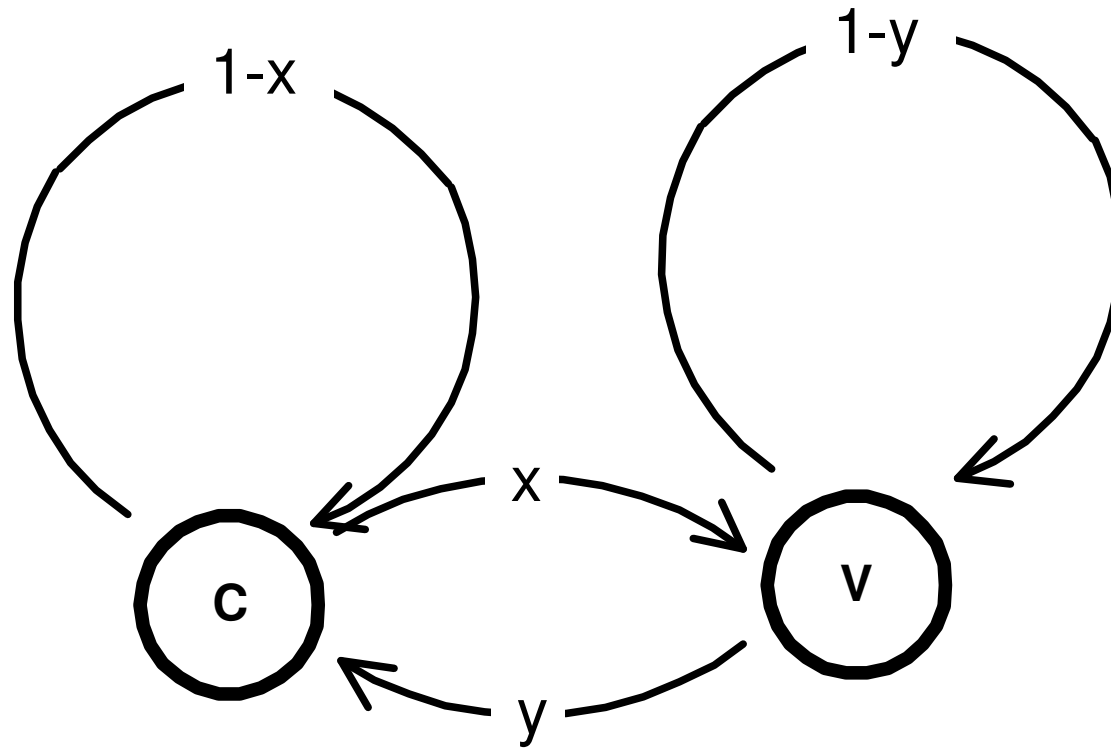
A note to the initiated

I am intentionally **leaving out** an important detail: **how to treat the word-final boundary ‘#’**. This is crucial to get the system to work correctly, but since I am assuming that this kind of model is not familiar to most of you, I will ignore this important point. If you *do* want to deal with it, you would modify the preceding model to become this:

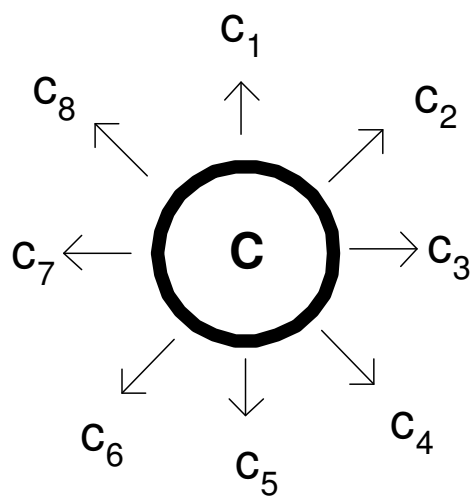
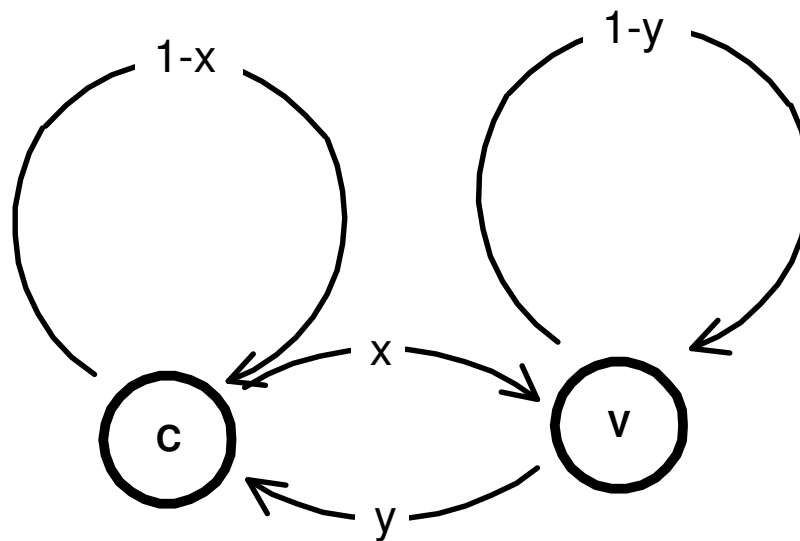


where State 1 generates all symbols *except* # (as before) and State 2 generates '#' and then the machine stops.

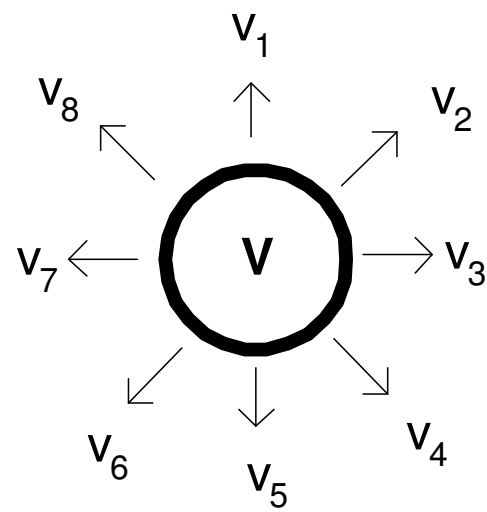
Much more interesting model:



For state transitions; and the same model for emissions: both states emit all of the symbols, but with different probabilities....



$$\sum_i c_i = 1$$



$$\sum_i v_i = 1$$

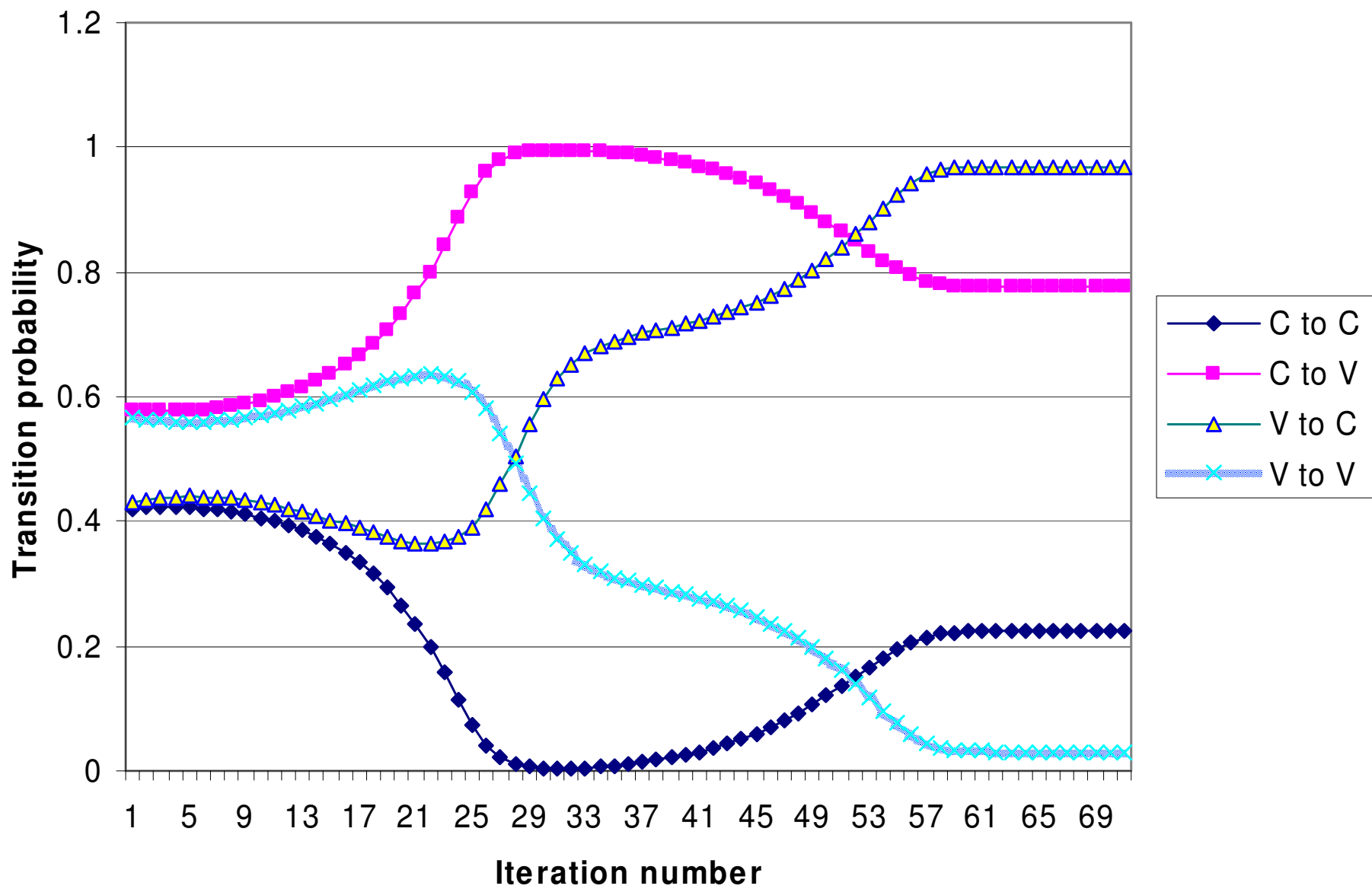
The question is...

- How could we obtain the *best* probabilities for p , q , and all of the emission probabilities for the two states?
- [Bear in mind: each state generates *all* of the symbols. The only way to ensure that a state does *not* generate a symbol is to assign a zero probability that the emission of the symbol in that state.]

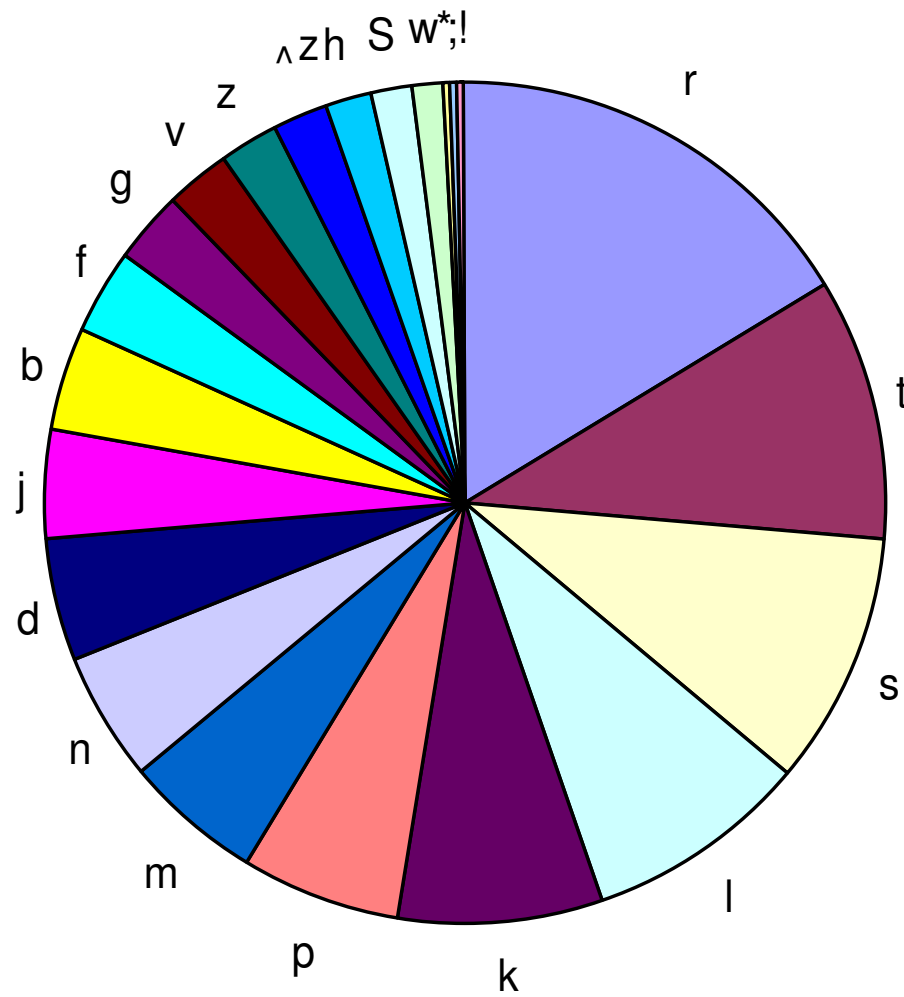
Results for 2 State HMM

- Separates Cs and Vs

Evolution of learning of 2 categories (3000 words)

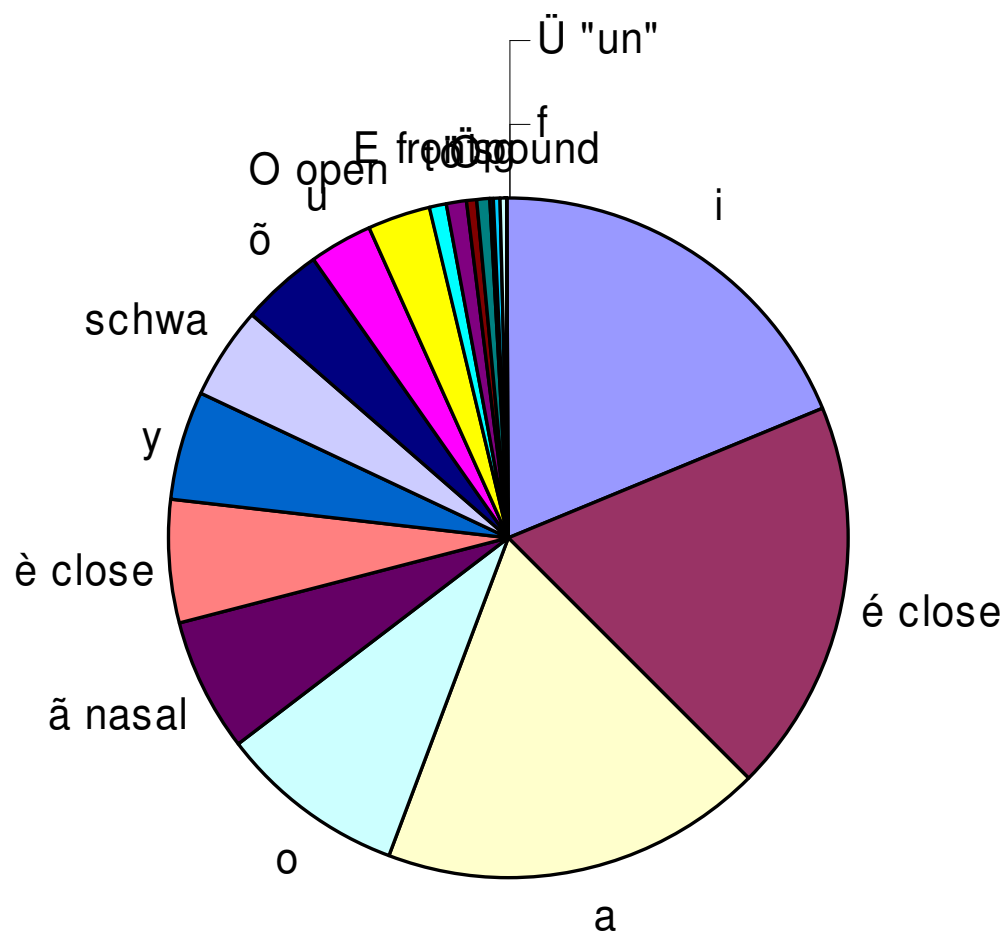


Members of Category 1 ("C")



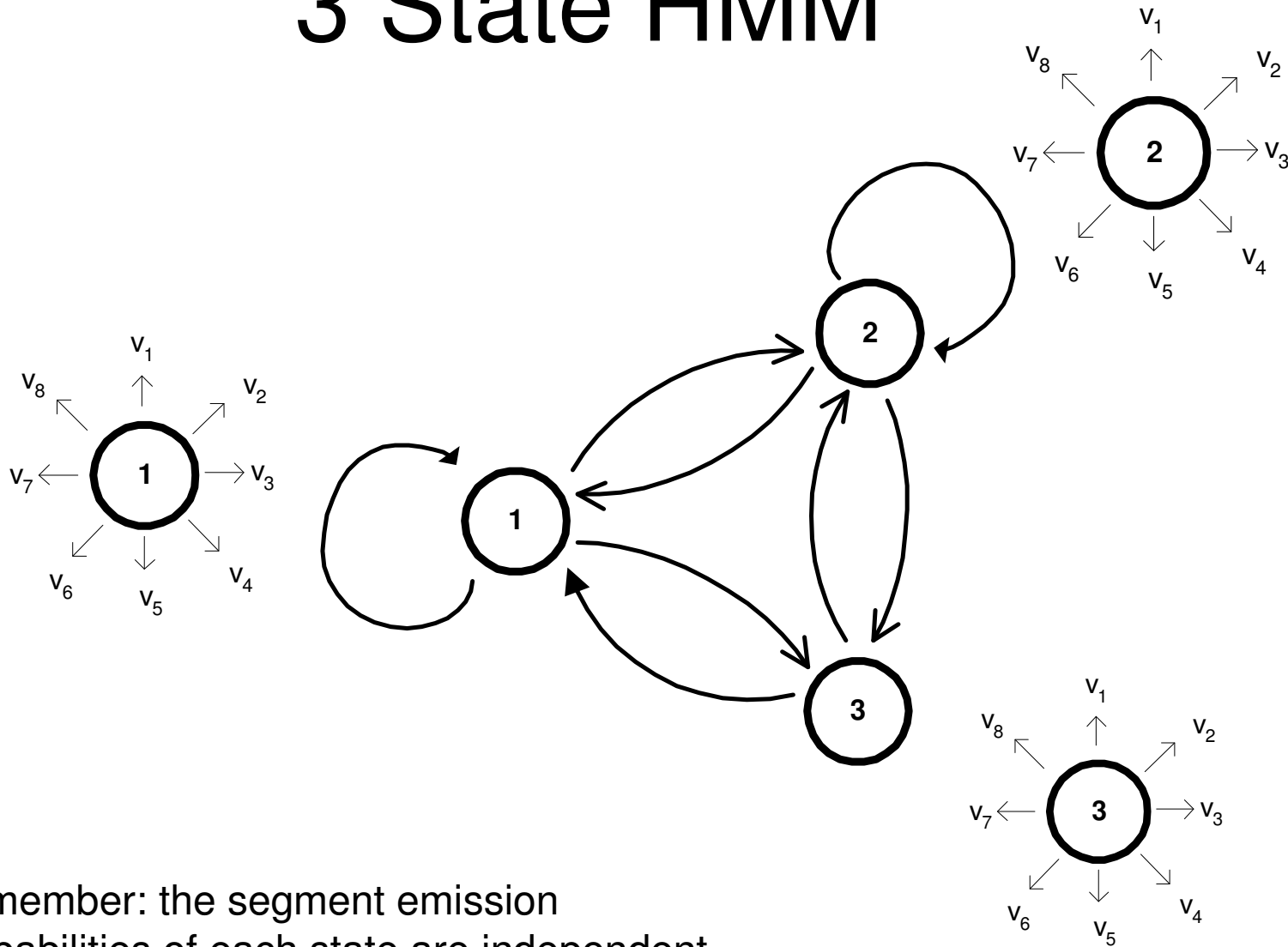
- r
- t
- s
- l
- k
- p
- m
- n
- d
- j
- b
- f
- g
- v
- z
- ^
- zh
- S
- w
- *

Members of Category 2 ("V")



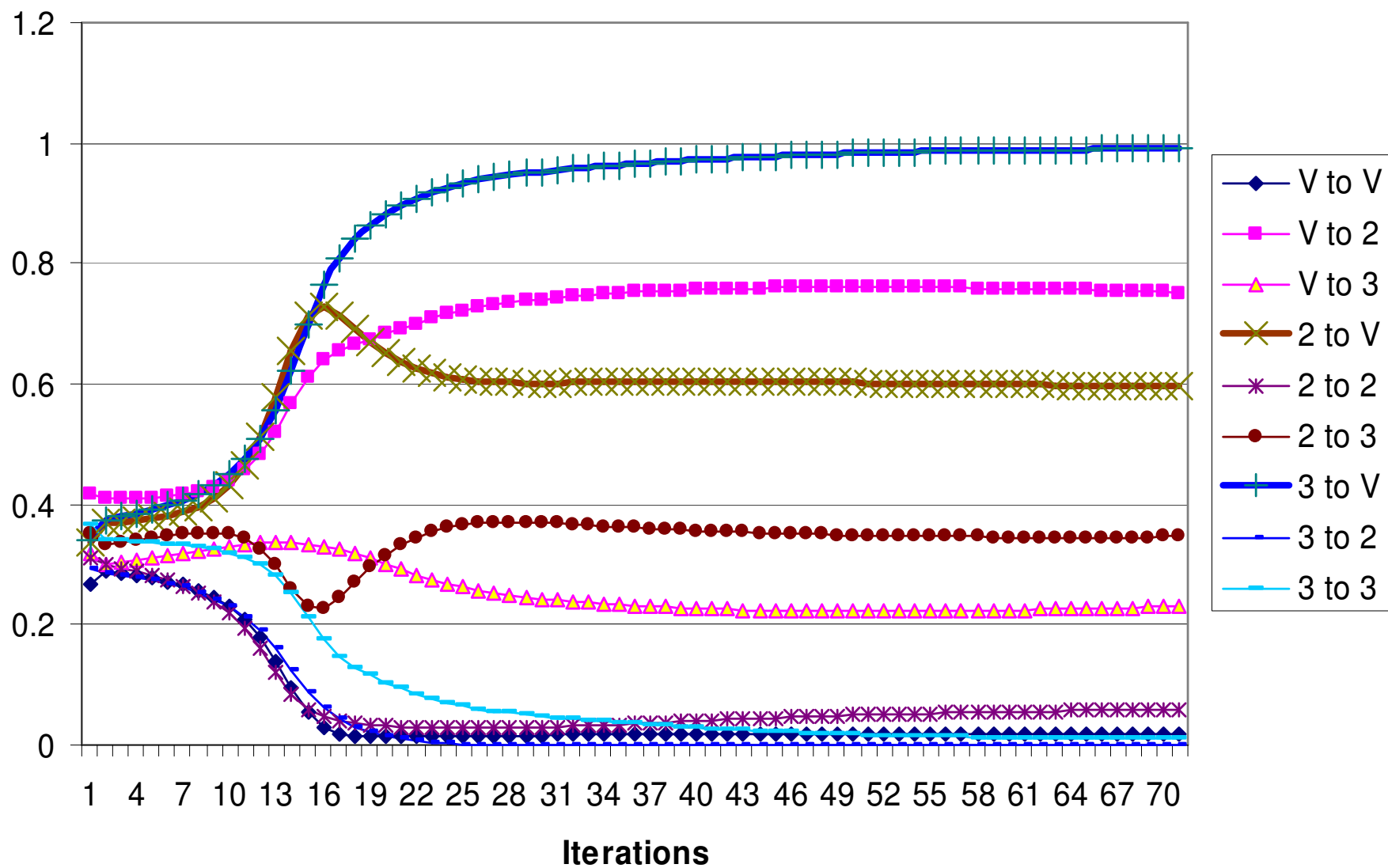
- i
- é close
- a
- o
- ã nasal
- è close
- y
- schwa
- õ
- u
- O open
- t
- ö
- Ö
- s
- E front round
- p
- g
- Ü "un"
- f

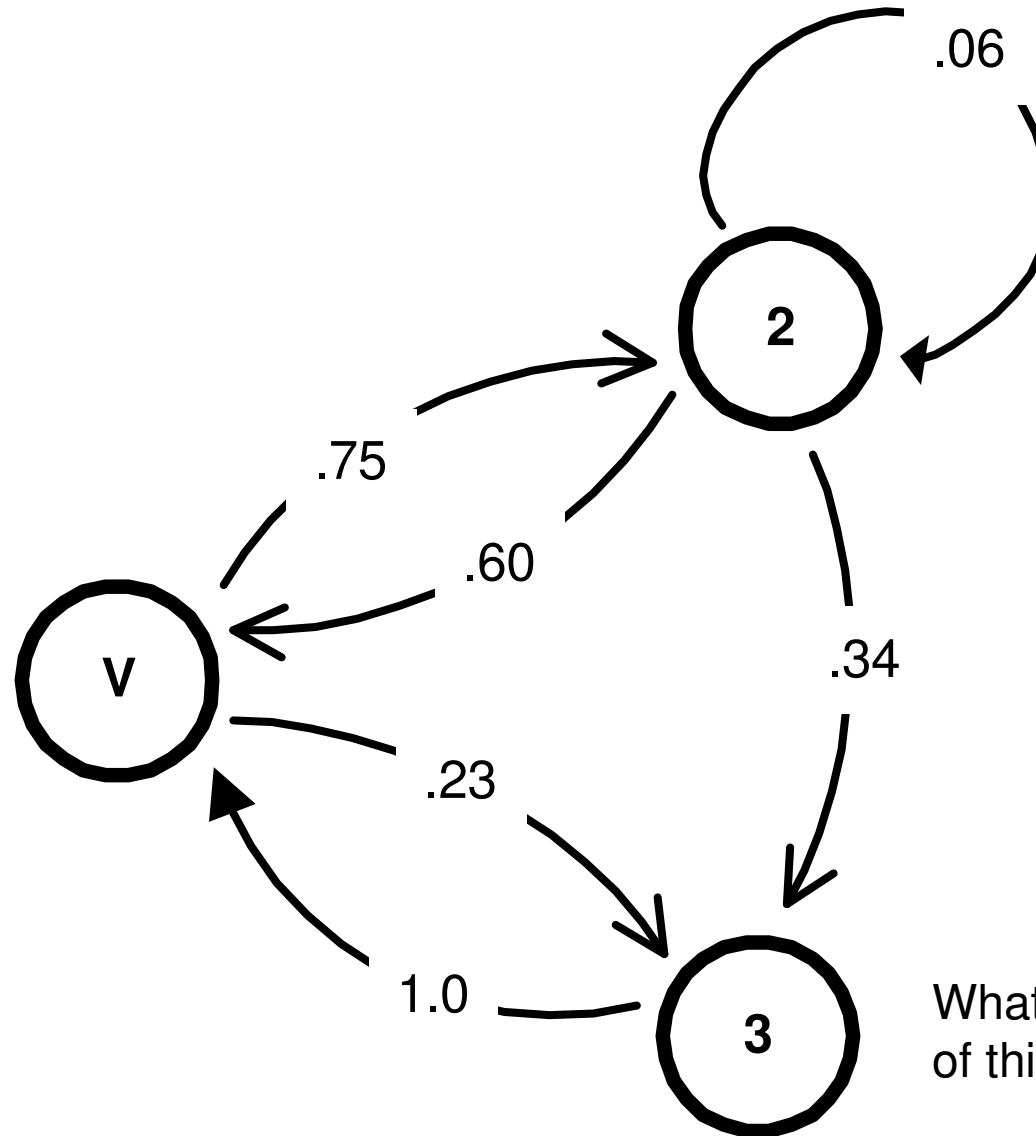
3 State HMM



Remember: the segment emission probabilities of each state are independent.

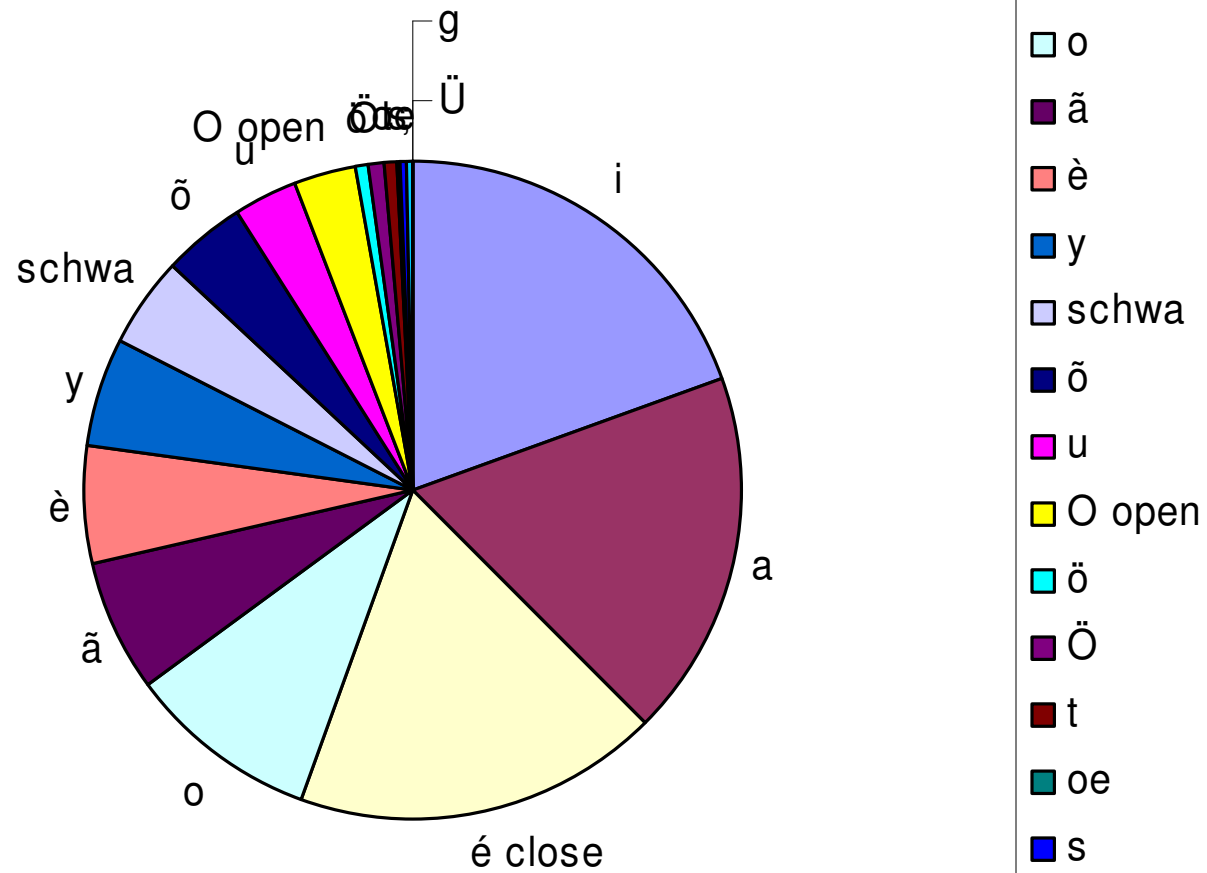
Learning of 3 state system



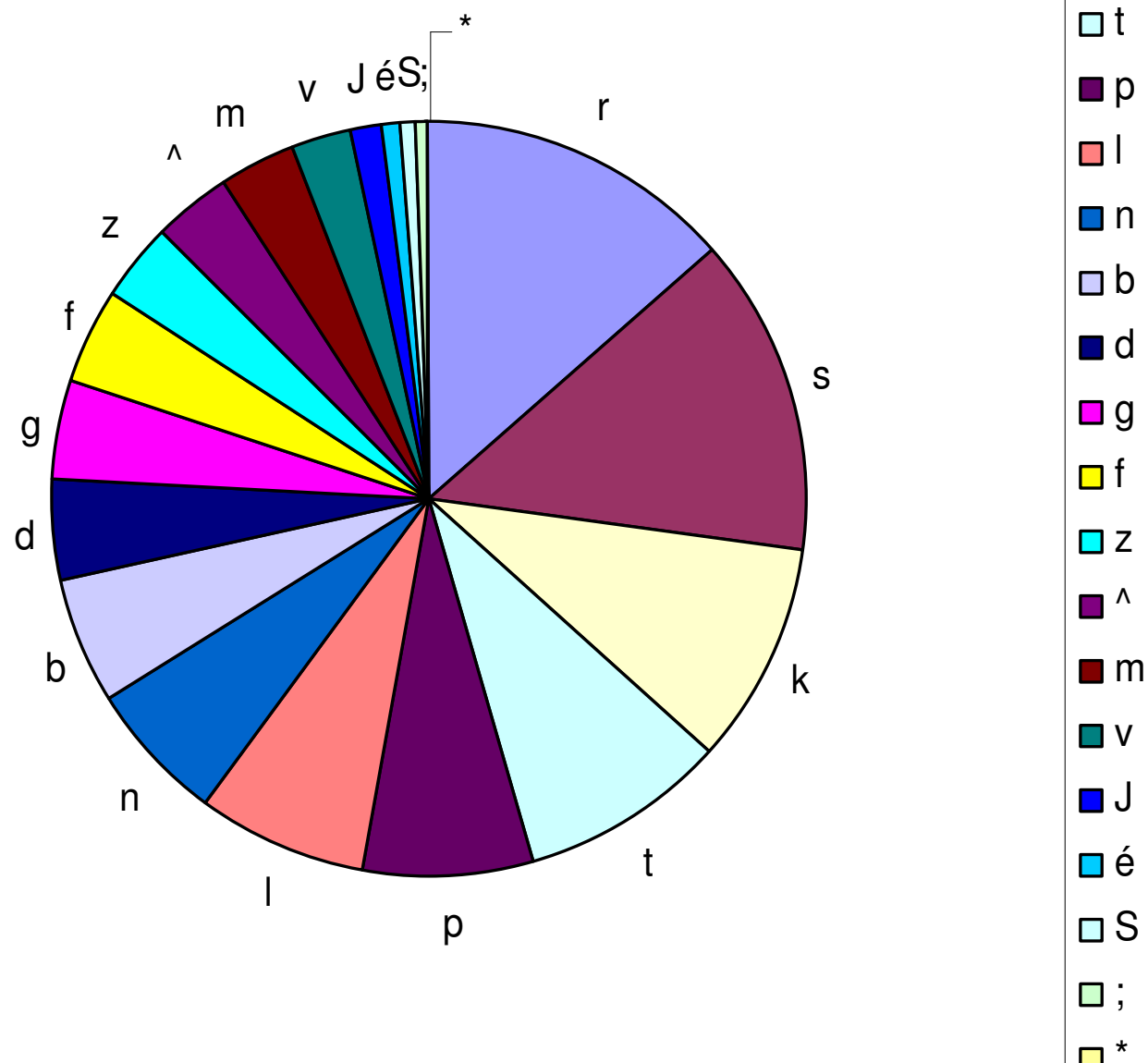


What is the “function” of this state?

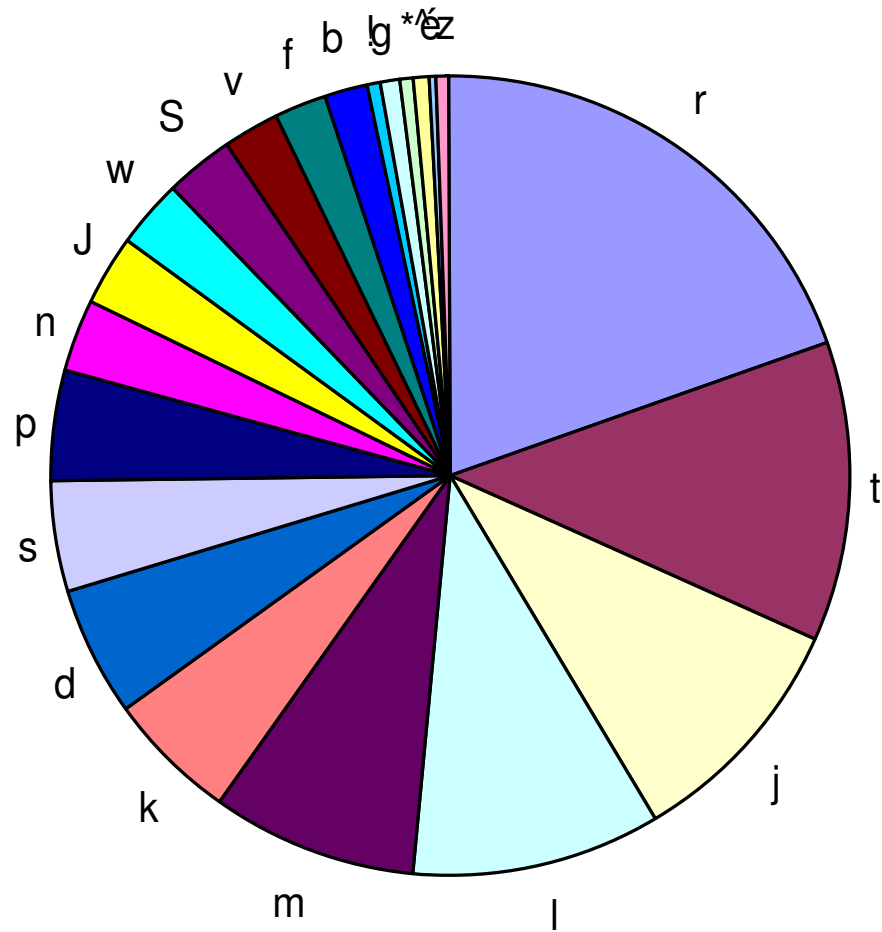
Members of Category 1 ("V")



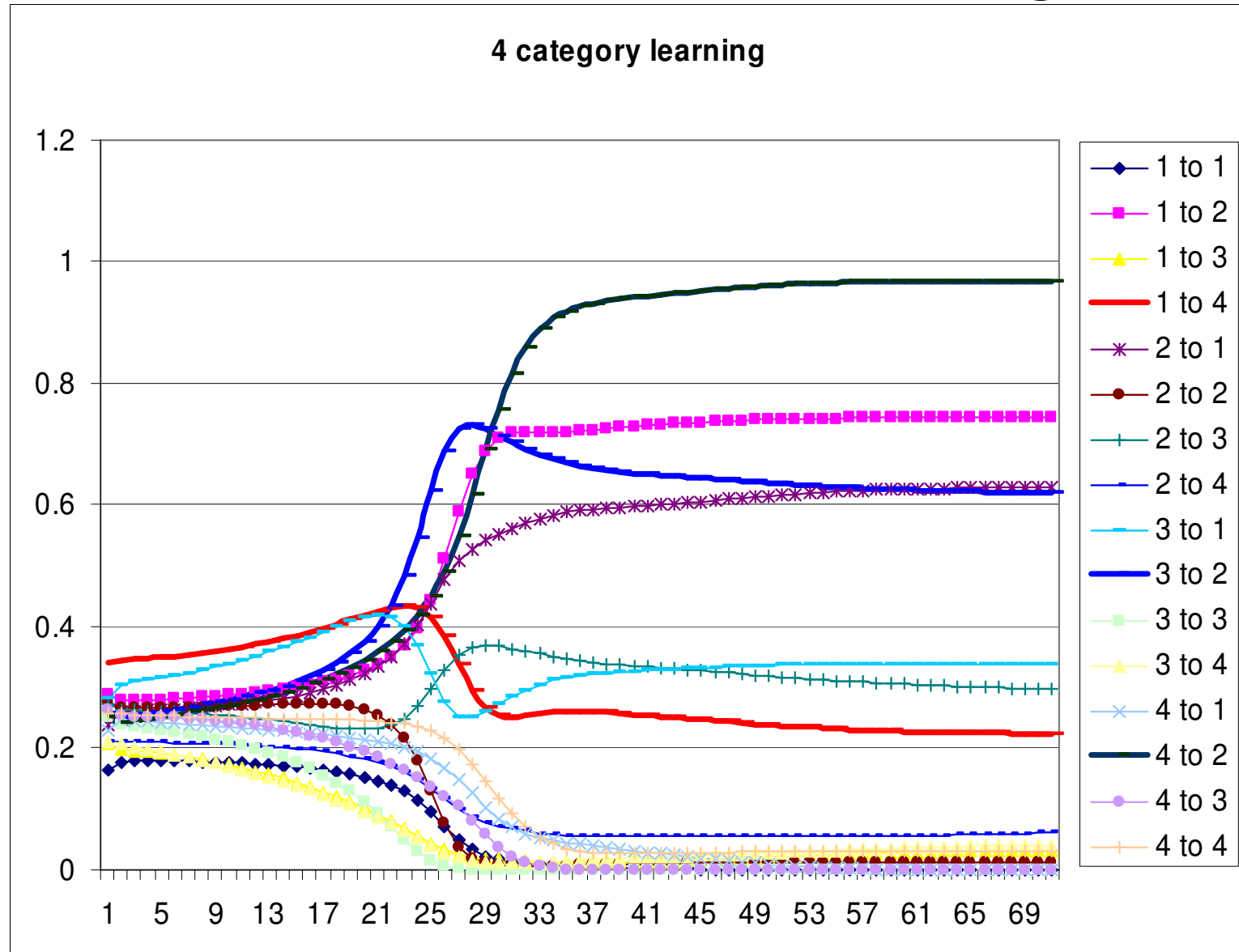
Members of Category 2

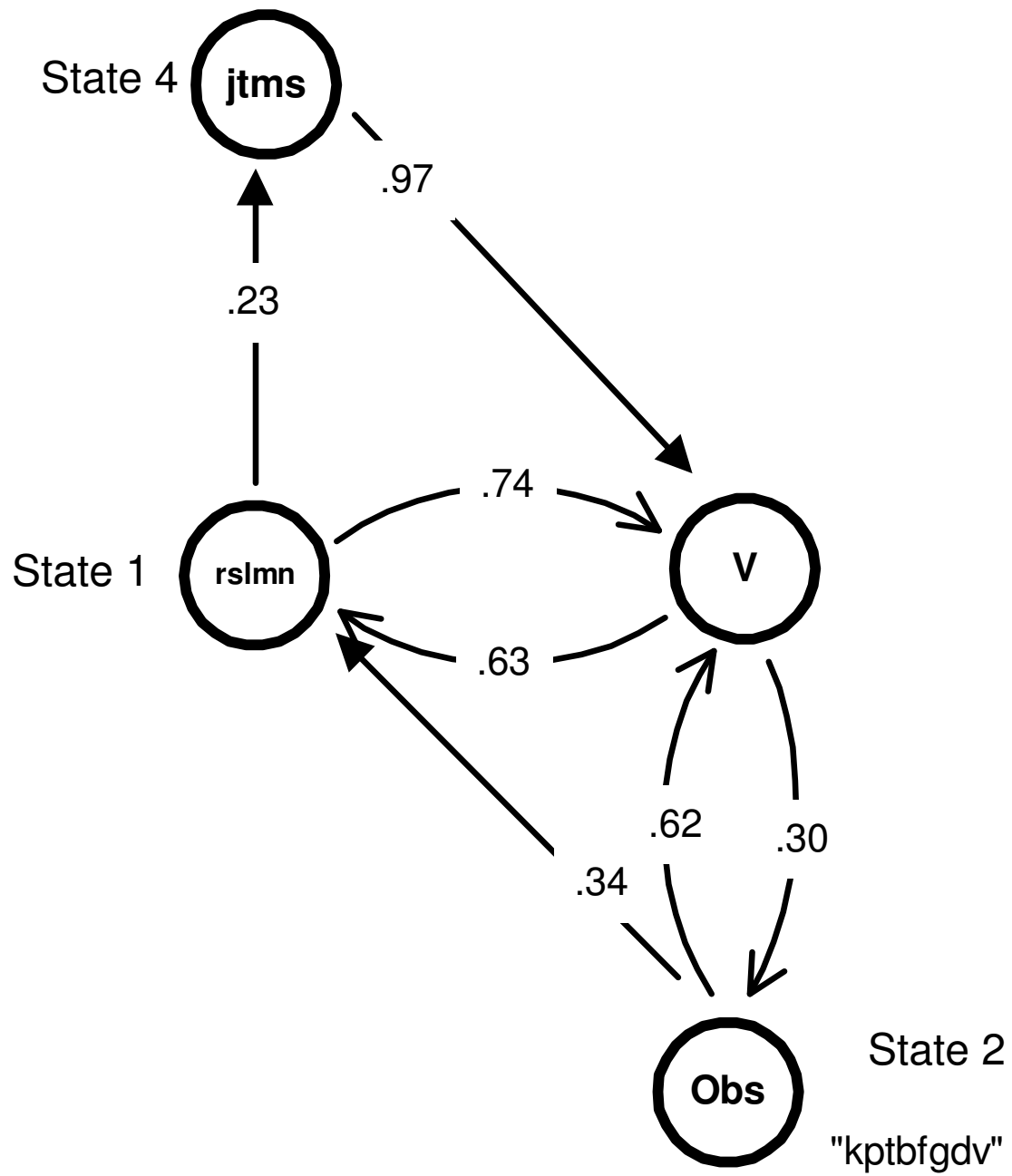


Members of Category 3

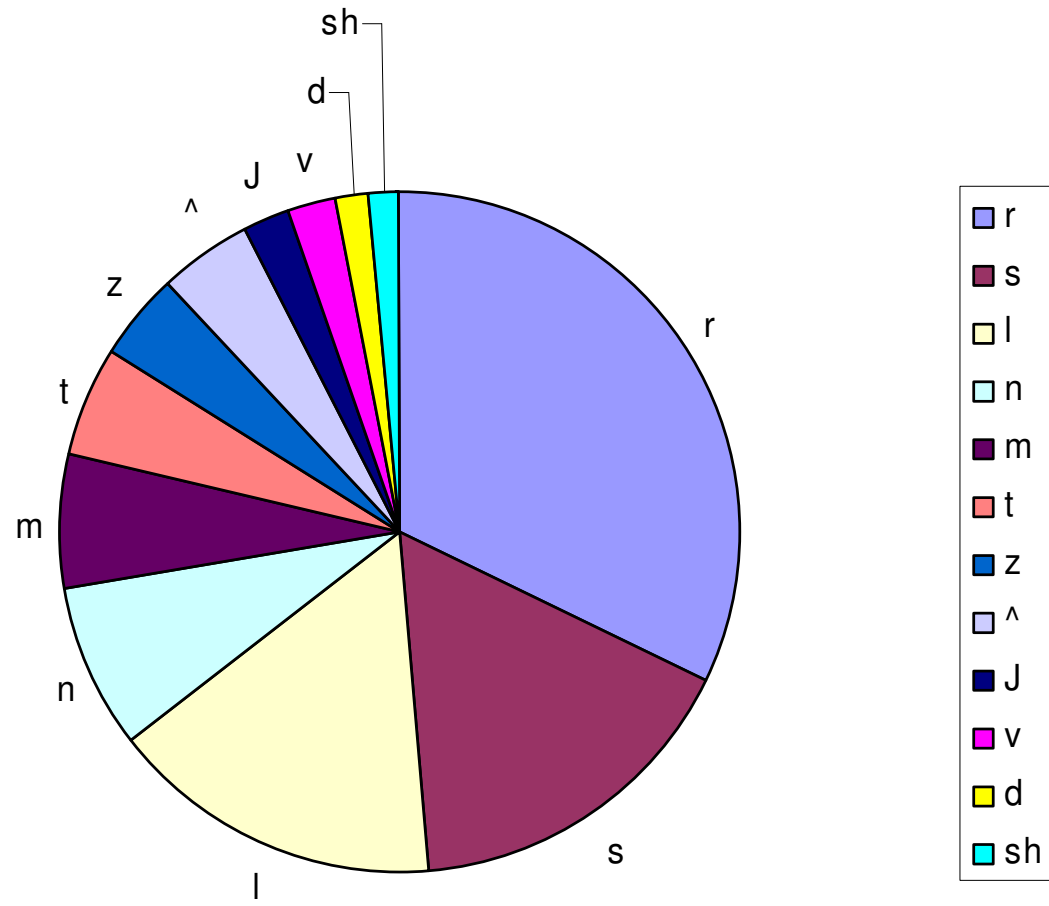


4 State HMM learning

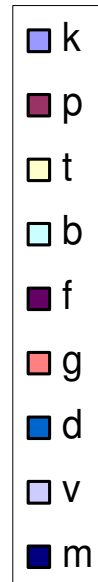
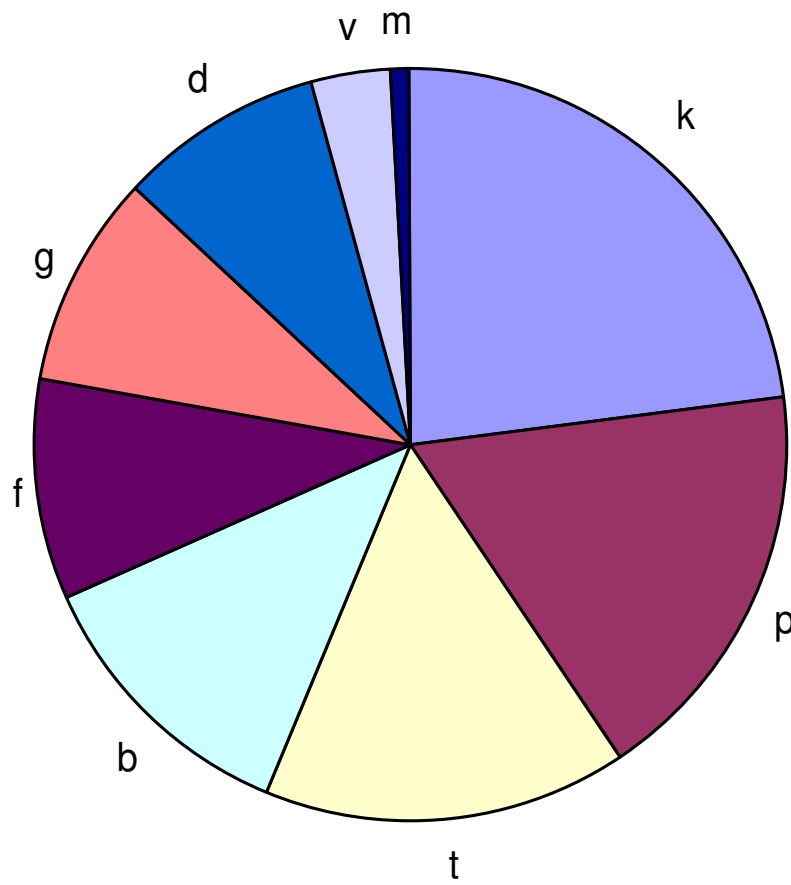




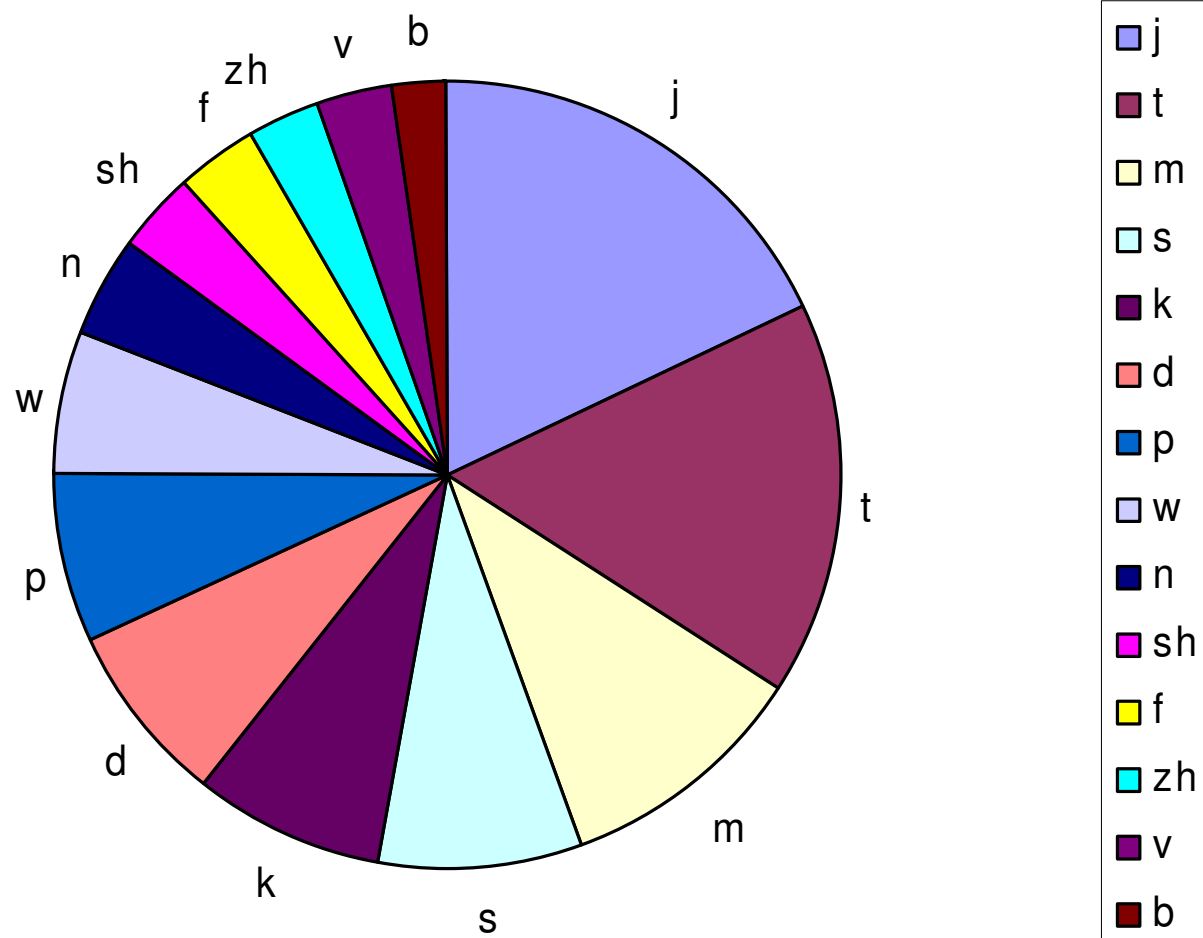
Category 1



Category 3



Category 4



Conclusion

- We have talked about *methods* instead of *theory*. Why?
- It is curious (isn't it?) that we tend to talk more about the development of phonological theory than about the development of phonological methods.
- We need both.

Phonological methods

- We need phonological methods when we want to use phonology to serve ends other than itself: e.g., to study variation.
- In addition: I am interested in these methods as tools to develop the capability for automatic learning of phonology.

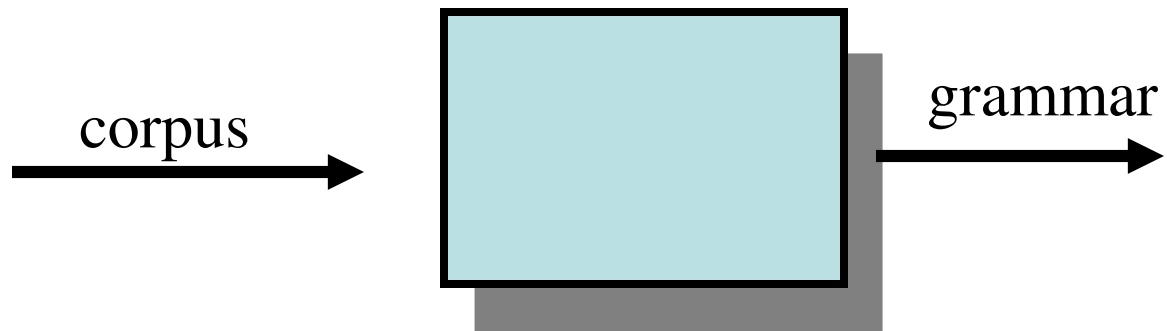
Automatic learning of grammar

- A return to the Harrisian goal of developing methods to project grammar from data.

End

Linguistic theory...

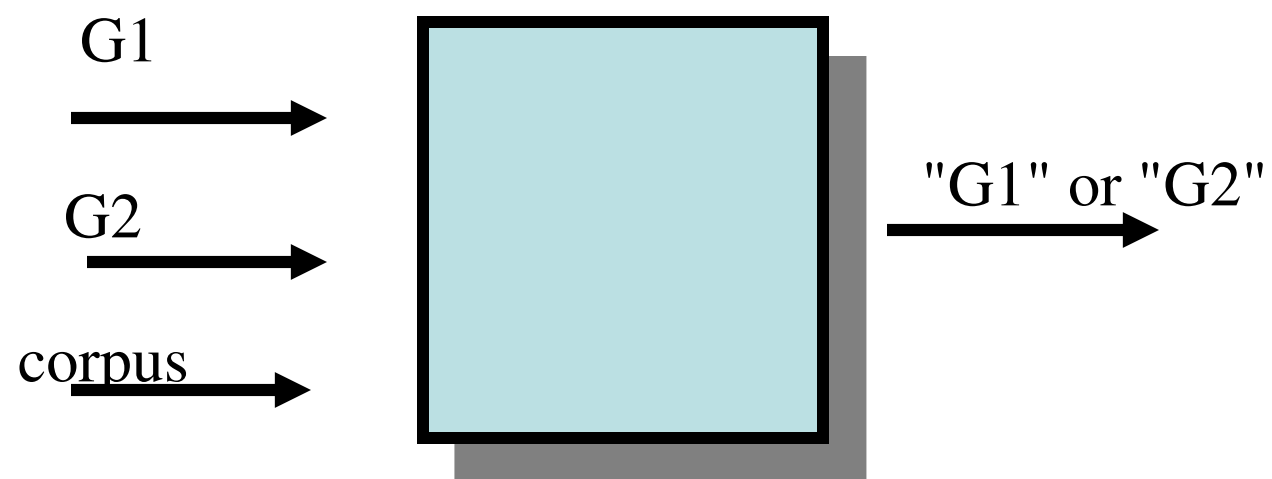
The strongest requirement that could be placed on the relation between a theory of linguistic structure and particular grammars is that the theory must provide a practical and mechanical method for actually constructing the grammar, given a corpus of utterances. Let us say that such a theory provides us with a *discovery procedure*.



- A weaker requirement would be that the theory must provide a practical and mechanical method for determining whether or not a grammar proposed for a given corpus is, in fact, the best grammar of the language from which the corpus is drawn (a *decision procedure*).



- An even weaker requirement would be that given a corpus and given two proposed grammars G_1 and G_2 , the theory must tell us which is the better grammar....an *evaluation procedure*.



The point of view adopted here is that it is unreasonable to demand of linguistic theory that it provide anything more than a practical evaluation procedure for grammars. That is, we adopt the weakest of the three positions described above...

I think that it is very questionable that this goal is attainable in any interesting way, and I suspect that any attempt to meet it will lead into a maze of more and more elaborate and complex analytic procedures that will fail to provide answers for many important questions about the nature of linguistic structure. I believe that *by lowering our*

sights to the more modest goal of developing an evaluation procedure for grammars we can focus attention more clearly on truly crucial problems...The correctness of this judgment can only be determined by the actual development and comparison of theories of these various sorts.

Notice, however, that the weakest of these three requirements is still strong enough to guarantee significance for a theory that meets it. **There are few areas of science in which one would seriously consider the possibility of developing a general, practical, mechanical method for choosing among several theories, each compatible with the available data.** (Noam Chomsky, *Syntactic Structures* (1957))