# Psych 229:
# Language Acquisition

Lecture 5
Statistics & Words

---

## Saffran, Aslin, & Newport 1996

*[slide text in small print]*

What does experience-independent mean (as opposed to experience-dependent)?

the task at hand

---

## Saffran, Aslin, & Newport 1996

transitional probability
= conditional probability

$$Y|X = \frac{\text{frequency of XY}}{\text{frequency of X}}$$

Why do they need the light, too?

---

## Saffran, Aslin, & Newport 1996

bidakupadotigolabubidaku

word divisions by transitional probability

p(ku pa) = 0.33

p(bi da) = 1.0

---

## Saffran, Aslin, & Newport 1996

**Table 1.** Mean time spent listening to the familiar and novel stimuli for experiment 1 (words versus nonwords) and experiment 2 (words versus part-words) and significance tests comparing the listening times.

| Experiment | Mean listening times (s) | | Matched-pairs t test |
|---|---|---|---|
| | Familiar items | Novel items | |
| 1 | 7.97 (SE = 0.41) | 8.85 (SE = 0.45) | t(23) = 2.3, P < 0.04 |
| 2 | 6.77 (SE = 0.44) | 7.60 (SE = 0.42) | t(23) = 2.4, P < 0.03 |

novelty preference

---

## Saffran, Aslin, & Newport 1996

So this isn't about learning *only* from transitional probabilities…

And this statistical tracking ability may apply only to language data…

Innate knowledge in the form of biases on learning, rather than explicit knowledge.

## Discussion Questions

How does statistical learning fit in with the idea of a mental grammar? What about with the idea of innateness? Can experience-independent mechanisms ensure learning by themselves in some situation?

Transitional probability: how does this fit into experience-dependent and experience-independent learning mechanisms?

## Gambell & Yang 2006: Computational model of word segmentation

**Survey of infant strategies (use at 8 months [before word meaning])**
Possible strategy: learn from isolated words

Data: 9% of mother-to-child speech is isolated words

Problem: How does a child recognize an isolated word as such?

length won't work: "I-see" vs. "spaghetti"

Possible strategy: statistical properties like transitional probability between syllables

word boundaries postulated at local minima

pre tty ba by          p(tty-->ba) < p(pre-->tty), p(ba-->by)

Question: How well does this fare on real data sets (not artificial stimuli)?

## Gambell & Yang 2006: Computational model of word segmentation

**Survey of infant strategies (use at 8 months [before word meaning])**
Possible strategy: Metrical segmentation strategy

Children treat stressed syllable as beginning of word

- 90% of English content words are stress-initial

Problem: Stress systems differ from language to language

- the child would need to know that words are stress initial

…but to do that, the child needs words *first*

Possible strategy: phonotactic constraints (sequences of consonant clusters that go together, e.g. **str** vs. **\*stl** in English); language-specific

- Infants seem to know these by 9 months

- posit boundary at improper sequence break: **stl** --> **st l** (fir**st l**ight)

Problem: May just be syllable boundary (re**stl**ess)

## Gambell & Yang 2006: Computational model of word segmentation

**Survey of infant strategies (use at 8 months [before word meaning])**
Possible strategy: Memory

Use previous stored words (sound forms, not meanings) to recognize new words

- if child knows *new*, then can recognize *one* in *thatsanewone*

Problem: Needs to know words before can use this

**A good point:** "It seems…only language-independent strategies can set word segmentation in motion before the establishment and application of language-specific strategies"

## Gambell & Yang 2006: Computational model of word segmentation

**Computational model goal**
- psychologically plausible learning algorithm

- real data

Another good point: it's good if the information is in the data, but we also need to know how children could use it

**On psychological plausibility**

## Gambell & Yang 2006: Computational model of word segmentation

**what to evaluate**

**where to get the data**

## Gambell & Yang 2006: Computational model of word segmentation

**modeling statistical learning (TPs)**

The modeling of statistical learning is straightforward, though it may be useful to make the details of our implementation clear. The model consists of two stages: training and testing. During the training stage, the learner gathers transitional probabilities over adjacent syllables in the learning data. The testing stage does not start until the entire learning data has been processed, and statistical learning is applied to the same data used in the training stage.

Another technical detail also needs to be spelled out: the TPs are gathered without stress information. That is, when counting syllable frequencies, the learner does not distinguish, say, a stressed syllable /ba/ from among the unstressed one."

That is, there is a word boundary AB and CD if if TP(A→B) >TP(B→C) < TP(C→D). The conjectured word boundaries are then compared against the target segmentation. Scoring is done for each utterance, using the definition of precision and recall in (1)

**results**

Modeling shows that the statistical learning (Saffran et al., 1996) does not reliably segment words such as those in child-directed English. Specifically, precision is 41.6%, recall is 23.3%. In other words, about 60% of words postulated by the statistical learner are not English words, and almost 80% of actual English words are not extracted. This is so even under favorable learning conditions:

- the child has syllabified the speech perfectly,
- the child has neutralized the effect of stress among the variants of syllables, which reduces the sparse data problem,
- and the data for segmentation is the same as the data used in training, which eliminates the sparse data problem

## Gambell & Yang 2006: Computational model of word segmentation

**What happened?**

We were surprised by the low level of performance. Upon close examination of the learning data, however, it is not difficult to understand the reason. A necessary condition on the use of TP local minima to extract words is that words must consist of multiple syllables. If the target sequence of segmentation contains only monosyllabic words, it is clear that statistical learning will fail. A sequence of monosyllabic words require a word boundary after each syllable; a statistical learner, on the other hand, will only place a word boundary between two sequences of syllables for which the TPs within are higher than that in the middle. Indeed, in the artificial language learning experiment of Saffran et al. (1996) and much subsequent work, the pseudowords are uniformly three syllables long. However, the case of child-directed English is quite different. The fact that the learning data consists of 226,178 words but only 263,660 syllables suggests that the overwhelming majority of word tokens are monosyllabic. More specifically, a monosyllabic word is followed by another monosyllabic word 85% of time. As long as this is the case, statistical learning cannot work.