

Learning English Metrical Phonology: Beyond Simple Probability

Lisa Pearl
 University of California, Irvine
 Sept 4, 2008
 GALANA 3

Human Language Learning

Theoretical work:
object of acquisition


(x	x)	x
H	L	H
EM	pha	sis

Human Language Learning

Theoretical work:
object of acquisition

(x	x)	x
H	L	H
EM	pha	sis

Experimental work:
time course of acquisition




Human Language Learning

Theoretical work:
object of acquisition

(x	x)	x
H	L	H
EM	pha	sis

Experimental work:
time course of acquisition



↗

mechanism of acquisition
given the boundary conditions provided by
(a) linguistic representation
(b) the trajectory of learning

Complex Linguistic Systems

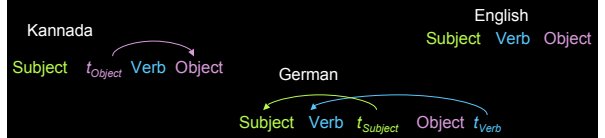
What is the generative system that creates the observed (structured) data of language (ex: **syntax**, metrical phonology)?

Observable data: **word order** Subject Verb Object

Complex Linguistic Systems

What is the generative system that creates the observed (structured) data of language (ex: **syntax**, metrical phonology)?

Observable data: **word order** Subject Verb Object



Complex Linguistic Systems

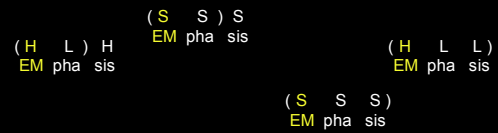
What is the generative system that creates the observed (structured) data of language (ex: syntax, **metrical phonology**)?

Observable data: **stress contour** EMphasis

Complex Linguistic Systems

What is the generative system that creates the observed (structured) data of language (ex: syntax, **metrical phonology**)?

Observable data: **stress contour** EMphasis



Road Map

Complex linguistic systems

- General problems
- Parametric systems
- Parametric metrical phonology
- Case study: English metrical phonology

Learnability of complex linguistic systems

- General learnability framework
- Previous learning successes: biased learners
- Unbiased probabilistic learning
- Where the problem lies: tricky data

Where next? Implications & Extensions

Road Map

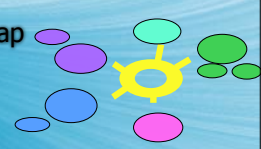
Complex linguistic systems

- General problems
- Parametric systems
- Parametric metrical phonology
- Case study: English metrical phonology

Learnability of complex linguistic systems


- General learnability framework
- Previous learning successes: biased learners
- Unbiased probabilistic learning
- Where the problem lies: tricky data

Where next? Implications & Extensions




General Problems with Learning Complex Linguistic Systems

What children encounter: the output of the generative linguistic system

EMphasis 

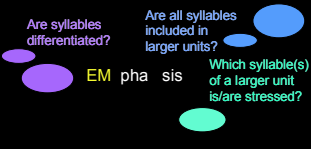
General Problems with Learning Complex Linguistic Systems

What children encounter: the output of the generative linguistic system

EMphasis 

What children must learn: the components of the system that combine to generate this observable output

- Are syllables differentiated?
- Are all syllables included in larger units?
- Which syllable(s) of a larger unit is/are stressed?



General Problems with Learning Complex Linguistic Systems

What children encounter: the output of the generative linguistic system

EMphasis



What children must learn: the components of the system that combine to generate this observable output

Are syllables differentiated?

Are all syllables included in larger units?

Which syllable(s) of a larger unit is/are stressed?

Why this is tricky:

There is often a non-transparent relationship between the observable form of the data and the underlying system that produced it. *Hard to know what parameters of variation to consider.*

(H L) H
EM pha sis

(S S S)
EM pha sis

The Hypothesis Space

Hypothesis for a language consists of a combination of generalizations about that language (grammar). But this leads to a theoretically infinite hypothesis space.

The Hypothesis Space

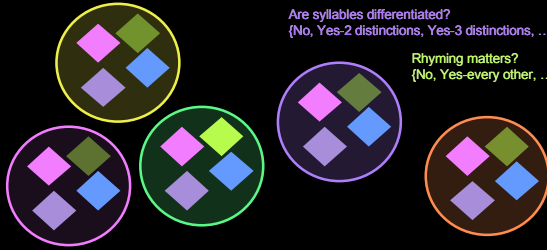
Hypothesis for a language consists of a combination of generalizations about that language (grammar). But this leads to a theoretically infinite hypothesis space.

Which syllable of a larger unit is stressed?
{Leftmost, Rightmost, Second from Left, ...}

Are all syllables included?
{Yes, No-not leftmost, No-not rightmost, ...}

Are syllables differentiated?
{No, Yes-2 distinctions, Yes-3 distinctions, ...}

Rhyming matters?
{No, Yes-every other, ...}



The Hypothesis Space

Hypothesis for a language consists of a combination of generalizations about that language (grammar). But this leads to a theoretically infinite hypothesis space.

Which syllable of a larger unit is stressed?
{Leftmost, Rightmost, ~~Second from Left, ...~~}

Are all syllables included?
{Yes, No-not leftmost, No-not rightmost, ~~...~~}

Are syllables differentiated?
{No, Yes-2 distinctions, Yes-3 distinctions, ~~...~~}

~~Rhyming matters?~~
{~~No, Yes-every other, ...~~}

Observation:
Languages only differ in constrained ways from each other. Not all generalizations are possible.



The Hypothesis Space

Hypothesis for a language consists of a combination of generalizations about that language (grammar). But this leads to a theoretically infinite hypothesis space.

Observation: Languages only differ in constrained ways from each other. Not all generalizations are possible.

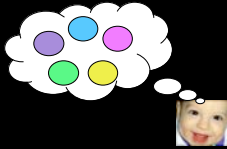
Idea: Bias on hypothesis space - children's hypotheses are constrained so they only consider generalizations that are possible in the world's languages.

Chomsky (1981), Halle & Vergnaud (1987), Tesar & Smolensky (2000)

Which syllable of a larger unit is stressed? {Leftmost, Rightmost}

Are all syllables included? {Yes, No-not leftmost, No-not rightmost}

Are syllables differentiated? {No, Yes-2 distinctions, Yes-3 distinctions}



Linguistic parameters = finite (if large) hypothesis space of possible grammars

Learning Parametric Linguistic Systems

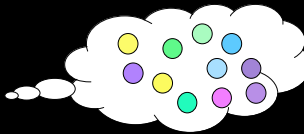
Linguistic parameters give the benefit of a finite hypothesis space. Still, the hypothesis space can be quite large.



(Clark 1994)

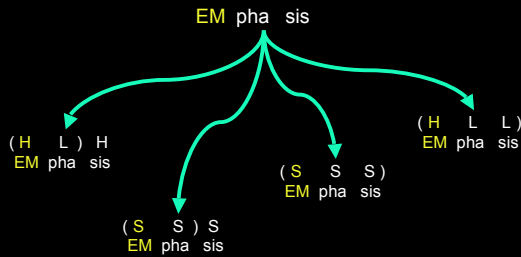
For example, assuming there are n binary parameters, there are 2^n core grammars to choose from.

Exponentially growing hypothesis space



Learning Parametric Linguistic Systems

Also, data are often *ambiguous* between competing hypotheses, since multiple grammars can account for the same data point.



Parametric Metrical Phonology

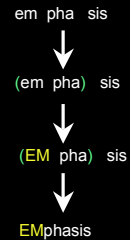
Metrical phonology: What tells you to put the **EM**phasis on a particular **SYL**lable

Process speakers use:
Basic input unit: syllables

Larger units formed: metrical feet
The way these are formed varies from language to language.

Stress assigned within metrical feet
The way this is done also varies from language to language.

Observable Data: stress contour of word **EM**phasis



Parametric Metrical Phonology

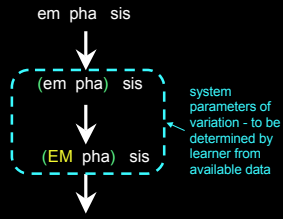
Metrical phonology:
What tells you to put the **EM**phasis on a particular **SYL**lable

Process speakers use:
Basic input unit: syllables

Larger units formed: metrical feet
The way these are formed varies from language to language.

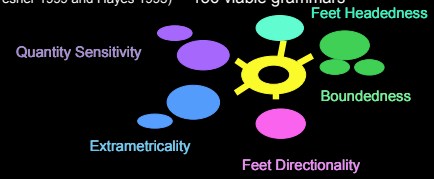
Stress assigned within metrical feet
The way this is done also varies from language to language.

Observable Data: stress contour of word **EM**phasis



Parametric Metrical Phonology

Metrical phonology system here: 5 main parameters, 4 sub-parameters
(adapted from Dresher 1999 and Hayes 1995) - 156 viable grammars

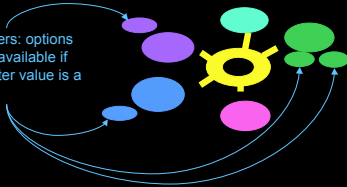


All combine to generate stress contour output

Parametric Metrical Phonology

Metrical phonology system here: 5 main parameters, 4 sub-parameters
(adapted from Dresher 1999 and Hayes 1995) - 156 viable grammars

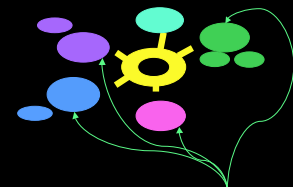
Sub-parameters: options that become available if main parameter value is a certain one



All combine to generate stress contour output

Parametric Metrical Phonology

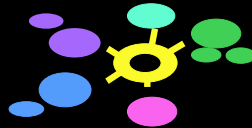
Metrical phonology system here: 5 main parameters, 4 sub-parameters
(adapted from Dresher 1999 and Hayes 1995) - 156 viable grammars



All combine to generate stress contour output

Generating a Stress Contour

Process speaker uses to generate stress contour

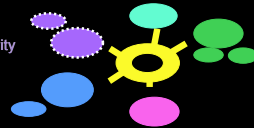


VC CV CVC
em pha sis

Generating a Stress Contour

Process speaker uses to generate stress contour

Quantity Sensitivity



Are syllables differentiated?

Yes - by rhyme.

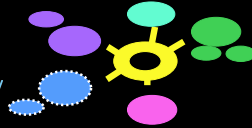
VC & VV syllables are Heavy, V syllables are Light.

H **L** **H**
VC CV CVC
em pha sis

Generating a Stress Contour

Process speaker uses to generate stress contour

Extrametricity



Are any syllables extrametrical?

Yes.

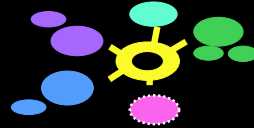
Rightmost syllable is not included in metrical foot.

(...)
H **L** **H**
VC CV CVC
em pha sis

Generating a Stress Contour

Process speaker uses to generate stress contour

Feet Directionality



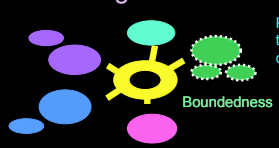
Which direction are feet constructed from?

From the right.


H **L** **H**
VC CV CVC
em pha sis

Generating a Stress Contour

Process speaker uses to generate stress contour



Boundedness



Are feet unrestricted in size?

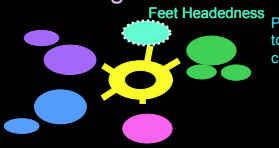
No.

2 syllables per foot.


(H	L)	H
VC	CV	CVC
em	pha	sis

Generating a Stress Contour

Process speaker uses to generate stress contour



Feet Headedness



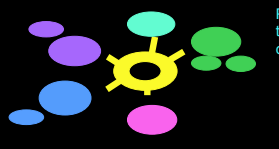

Which syllable of the foot is stressed?

Leftmost.

(H	L)	H
VC	CV	CVC
em	pha	sis


Generating a Stress Contour

Process speaker uses to generate stress contour

Learner's task: Figure out which parameter values were used to generate this contour.


(H	L)	H
VC	CV	CVC
EM	pha	sis



Case study: English metrical phonology

Estimate of child input: caretaker speech to children between the ages of 6 months and 2 years (CHILDES [Brent & Bernstein corpora]: MacWhinney 2000)


Total Words: 540505 Mean Length of Utterance: 3.5



Words parsed into syllables using the MRC Psycholinguistic database (Wilson, 1988) and assigned likely stress contours using the American English CALLHOME database of telephone conversation (Canavan et al., 1997)

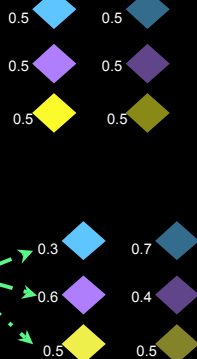
Road Map

- Complex linguistic systems
 - General problems
 - Parametric systems
 - Parametric metrical phonology
 - Case study: English metrical phonology
- Learnability of complex linguistic systems**
 - General learnability framework
 - Previous learning successes: biased learners
 - Unbiased probabilistic learning
 - Where the problem lies: tricky data
- Where next? Implications & Extensions




The learning framework: 3 components

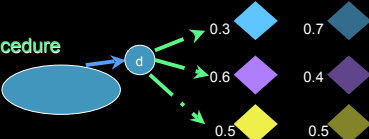
(1) Hypothesis space



(2) Data



(3) Update procedure

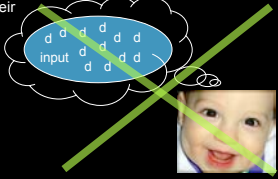


Key point for cognitive modeling: psychological plausibility

Any probabilistic update procedure that children are likely to use must, at the very least, be **incremental/online**.

Why? Humans (especially human children) don't have infinite memory.

Unlikely: human children can hold a whole corpus's worth of data in their minds for analysis later on



Key point for cognitive modeling: psychological plausibility

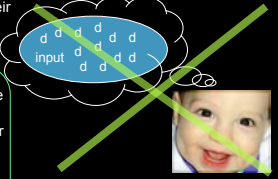
Any probabilistic update procedure that children are likely to use must, at the very least, be **incremental/online**.

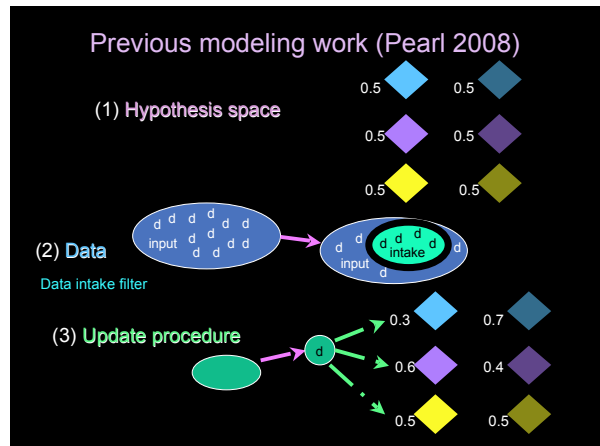
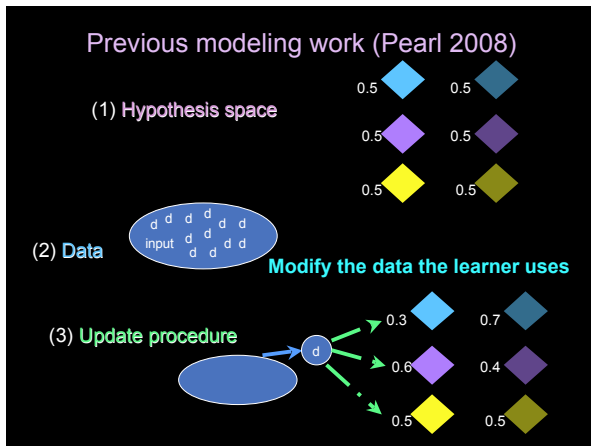
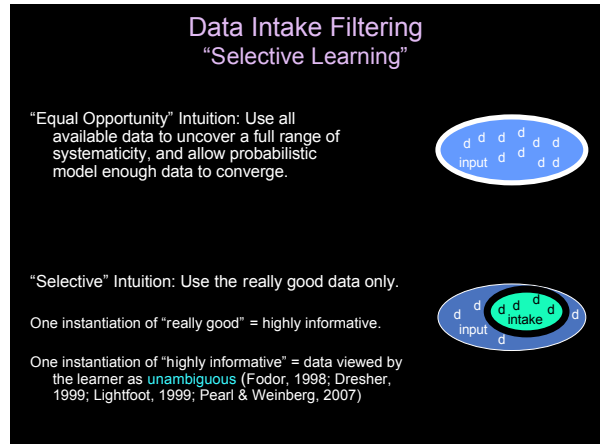
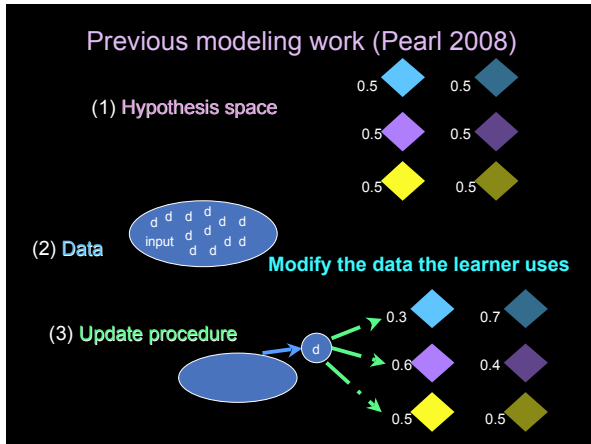
Why? Humans (especially human children) don't have infinite memory.

Unlikely: human children can hold a whole corpus's worth of data in their minds for analysis later on

Learning algorithms that operate over an entire data set do not have this property.
(ex: Foraker et al. 2007, Goldwater et al. 2007)

Desired: Learn from a single data point, or perhaps a small number of data points at most.





Biased learner, using only unambiguous data

Pearl (2008): Success is guaranteed as long as the parameters are learned in a particular order.

However...this requires the learner to identify unambiguous data and know/derive the appropriate parameter-setting order, which may not be trivial.

So...is this selective learning bias really necessary?
How well do unbiased learners do?

Two psychologically plausible probabilistic update procedures



Yang (2002)

Naïve Parameter Learner (**NParLearner**)

Probabilistic generation & testing of grammars. (incremental)
Hypothesis update: **Linear reward-penalty**

(Bush & Mosteller 1951)

Two psychologically plausible probabilistic update procedures



Yang (2002)

Naïve Parameter Learner (**NParLearner**)

Probabilistic generation & testing of grammars. (incremental)
Hypothesis update: **Linear reward-penalty**

(Bush & Mosteller 1951)



MAP Bayesian Learner (**BayesLearner**)

Probabilistic generation & testing of grammars. (incremental)
Hypothesis update: **Bayesian updating**

(Chew 1971: binomial distribution)

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

For each parameter, the learner associates a probability with each of the competing parameter values.

QI = 0.5	QS = 0.5
QSVCL = 0.5	QSVCH = 0.5
Em-Some = 0.5	Em-None = 0.5
Em-Left = 0.5	Em-Right = 0.5
Ft Dir Left = 0.5	Ft Dir Rt = 0.5
Bounded = 0.5	Unbounded = 0.5
Bounded-2 = 0.5	Bounded-3 = 0.5
Bounded-Syl = 0.5	Bounded-Mor = 0.5
Ft Hd Left = 0.5	Ft Hd Rt = 0.5

↑
Initially all are equiprobable

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

For each data point encountered, the learner probabilistically generates a grammar.

AFTERNOON

QI = 0.5	QS = 0.5
QSVCL = 0.5	QSVCH = 0.5
Em-Some = 0.5	Em-None = 0.5
Em-Left = 0.5	Em-Right = 0.5
Ft Dir Left = 0.5	Ft Dir Rt = 0.5
Bounded = 0.5	Unbounded = 0.5
Bounded-2 = 0.5	Bounded-3 = 0.5
Bounded-Syl = 0.5	Bounded-Mor = 0.5
Ft Hd Left = 0.5	Ft Hd Rt = 0.5

QI/QS?...if QS, QSVCL or QSVCH?
Em-None/Em-Some?...

QS, QSVCL, Em-None, Ft Dir Right,
Bounded, Bounded-2, Bounded-Syl, Ft Hd Right

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

The learner then uses this grammar to generate a stress contour for the observed data point.

AFTERNOON

If the generated stress contour matches the observed stress contour, all participating parameter values are rewarded.

QS, QSVCL, Em-None, Ft Dir Right, Bounded, Bounded-2, Bounded-Syl, Ft Hd Right	→	(L) (L) (H)
		VC CVC CVVC
		af ter NOON

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

The learner then uses this grammar to generate a stress contour for the observed data point.

AFTERNOON

If the generated stress contour does *not* match the observed stress contour, all participating parameter values are punished.

QS, QSVCL, Em-None, Ft Dir Left, Bounded, Bounded-2, Bounded-Syl, Ft Hd Right	→	(L) (L) (H)
		VC CVC CVVC
		af TER NOON

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

The learner then uses this grammar to generate a stress contour for the observed data point.

AFTERNOON

If the generated stress contour matches the observed stress contour, all participating parameter values are rewarded.

QS, QSVCL, Em-None, Ft Dir Right, Bounded, Bounded-2, Bounded-Syl, Ft Hd Right	→	(L) (L) (H)
		VC CVC CVVC
		af ter NOON

Match (success): reward all

QS, QSVCL, Em-None, Ft Dir Left, Bounded, Bounded-2, Bounded-Syl, Ft Hd Right	→	(L) (L) (H)
		VC CVC CVVC
		af TER NOON

Mismatch (failure): punish all

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities

NParLearner (Yang 2002): Linear Reward-Penalty

Learning rate γ :
small = small changes
large = large changes

Parameter values v1 vs. v2	
$p_{v1} = p_{v1} + \gamma(1 - p_{v1})$	$p_{v1} = (1 - \gamma)p_{v1}$
$p_{v2} = 1 - p_{v1}$	$p_{v2} = 1 - p_{v1}$
reward v1	punish v1

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities

NParLearner (Yang 2002): Linear Reward-Penalty

Learning rate γ :
small = small changes
large = large changes

Parameter values v1 vs. v2	
$p_{v1} = p_{v1} + \gamma(1 - p_{v1})$	$p_{v1} = (1 - \gamma)p_{v1}$
$p_{v2} = 1 - p_{v1}$	$p_{v2} = 1 - p_{v1}$
reward v1	punish v1

BayesLearner: Bayesian update of binomial distribution (Chew 1971)

Parameters α, β :

$\alpha = \beta$: initial bias at $p = 0.5$
 $\alpha, \beta < 1$: initial bias toward endpoints ($p = 0.0, 1.0$)

Parameter value v1	
$p_v = \frac{\alpha + 1 + \text{successes}}{\alpha + \beta + 2 + \text{total data seen}}$	
reward: success + 1	punish: success + 0

here: $\alpha = \beta = 0.5$

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities

After learning: expect probabilities of parameter values to converge near endpoints (above/below some threshold).

QI = 0.3	QS = 0.7
QSVCL = 0.6	QSVCH = 0.4
Em-Some = 0.1	Em-None = 0.9
...	...

Probabilistic learning for English

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities

After learning: expect probabilities of parameter values to converge near endpoints (above/below some threshold).

QI = 0.3	QS = 0.7
QSVCL = 0.6	QSVCH = 0.4
Em-Some = 0.1	Em-None = 0.9
...	...

Once set, a parameter value is always used during generation, since its probability is 1.0. Em-None = 1.0
(Em-Some = 0.0)

QI/QS?...if QS, QSVCL or QSVCH?
Em-None
...

→ QS, QSVCL, Em-None, Ft Dir Right,
Bounded, Bounded-2, Bounded-Syl, Ft Hd Right

Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%



Examples of incorrect target grammars

NParLearner:

Em-None, Ft Hd Left, Unb, Ft Dir Left, QI

QS, Em-None, QSVCH, Ft Dir Rt, Ft Hd Left, B-Mor, Bounded, Bounded-2

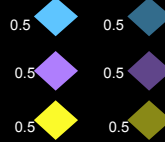
BayesLearner:

QS, Em-Some, Em-Right, QSVCH, Ft Hd Left, Ft Dir Rt, Unb

Bounded, B-Syl, QI, Ft Hd Left, Em-None, Ft Dir Left, B-2

The learning framework: where can we modify?

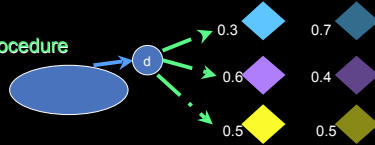
(1) Hypothesis space



(2) Data

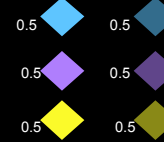


(3) Update procedure



The learning framework: where can we modify?

(1) Hypothesis space

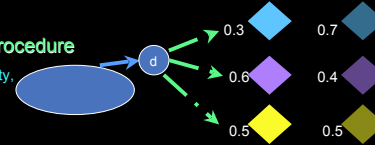


(2) Data



(3) Update procedure

Linear Reward-Penalty, Bayesian...?



Probabilistic learning for English: Modifications

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities

Batch-learning (for very small batch sizes): smooth out some of the irregularities in the data, better deal with complex systems (Yang 2002)

Implementation (Yang 2002):

Matching contour = increase parameter value's batch counter by 1
 Mismatching contour = decrease parameter value's batch counter by 1

Invoke update procedure (Linear Reward-Penalty or Bayesian Updating) **when batch limit b is reached.**

Probabilistic learning for English: Modifications

Probabilistic generation and testing of grammars (Yang 2002)

Update parameter value probabilities + Batch Learning

NParLearner (Yang 2002): Linear Reward-Penalty

Invoke when the batch counter for p_{v1} or p_{v2} equals b .

Parameter values $v1$ vs. $v2$

$$p_{v1} = p_{v1} + \gamma(1 - p_{v1}) \quad p_{v1} = (1 - \gamma)p_{v1}$$

$$p_{v2} = 1 - p_{v1} \quad p_{v2} = 1 - p_{v1}$$

reward $v1$ punish $v1$

BayesLearner: Bayesian update of binomial distribution (Chew 1971)

Invoke when the batch counter for p_{v1} or p_{v2} equals b .

Parameter value $v1$

$$p_{v1} = \frac{\alpha + 1 + \text{successes}}{\alpha + \beta + 2 + \text{total data seen}}$$

Note: total data seen + 1

reward: success + 1 punish: success + 0

Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%



Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%
NParLearner + Batch, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	0.8%
BayesLearner + Batch, $2 \leq b \leq 10$	1.0%



Probabilistic learning for English: Modifications

Probabilistic generation and testing of grammars (Yang 2002)

Learner hypothesis bias: metrical phonology relies in part on knowledge of rhythmical properties of the language

Human infants may already have knowledge of Ft Hd Left and QS.

Jusczyk, Cutler, & Redanz (1993): English 9-month olds prefer strong-weak stress bisyllables (trochaic) to weak-strong ones (iambic).



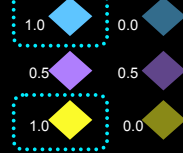
Turk, Jusczyk, & Gerken (1995): English infants are sensitive to the difference between long vowels and short vowels in syllables



Where else can we modify?

(1) Hypothesis space

Prior knowledge, biases:
QS, Ft Hd Left known...

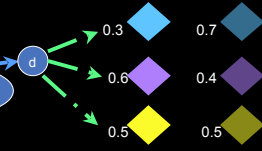


(2) Data



(3) Update procedure

Linear Reward-Penalty,
Bayesian, Batch...



Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%
NParLearner + Batch, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	0.8%
BayesLearner + Batch, $2 \leq b \leq 10$	1.0%



Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%
NParLearner + Batch, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	0.8%
BayesLearner + Batch, $2 \leq b \leq 10$	1.0%
NParLearner + Batch + Bias, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	5.0%
BayesLearner + Batch + Bias, $2 \leq b \leq 10$	1.0%



Probabilistic learning for English

Goal: Converge on English values after learning period is over

Learning Period Length: 1,666,667 words (based on estimates of words heard in a 6 month period, using Akhtar et al. (2004)).

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Model	Success rate (1000 runs)
NParLearner, $0.01 \leq \gamma \leq 0.05$	1.2%
BayesLearner	0.0%
NParLearner + Batch, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	0.8%
BayesLearner + Batch, $2 \leq b \leq 10$	1.0%
NParLearner + Batch + Bias, $0.01 \leq \gamma \leq 0.05, 2 \leq b \leq 10$	5.0%
BayesLearner + Batch + Bias, $2 \leq b \leq 10$	1.0%

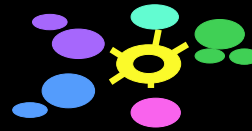


The best isn't so great

What gives?

Metrical phonology system here: 5 main parameters, 4 sub-parameters (adapted from Dresher 1999 and Hayes 1995)

156 viable grammars



English is *not* the optimal grammar

Adult English system values:

QS, QSVCH, Em-Some, Em-Right, Ft Dir Right, Bounded, Bounded-2, Bounded-Syllabic, Ft Hd Left

Of the 156 available grammars, English is ranked

52nd by token compatibility

56th by type compatibility

If prior knowledge of the hypothesis space is assumed (Ft Hd Left and QS), there are 60 available grammars.

English is ranked

18th by token compatibility

18th by type compatibility

Unbiased probabilistic learning is more likely to find the optimal grammar

English is compatible with 72.97% of the data by tokens, and 62.14% of the data by types.

The average compatibility of the grammars selected by unbiased probabilistic learning (using batch learning) was 73.56% of the data by tokens and 63.3% of the data by types.

Unbiased probabilistic learning is more likely to find the optimal grammar

English is compatible with **72.97%** of the data by tokens, and **62.14%** of the data by types.

The average compatibility of the grammars selected by unbiased probabilistic learning (using batch learning) was **73.56%** of the data by tokens and **63.3%** of the data by types.

Unbiased probabilistic learning works just fine - it's the English child-directed speech that's the problem!

Biased Children

The data actually lead an unbiased probabilistic learner to more optimal grammars than the English grammar.

Yet English children seem to learn the English grammar.

Conclusion: Children must have some additional bias that causes the sub-optimal English grammar to become the optimal grammar for this data set.

One idea: selective learning bias to heed only unambiguous data (Pearl 2008)



Road Map

Complex linguistic systems

- General problems
- Parametric systems
- Parametric metrical phonology
- Case study: English metrical phonology

Learnability of complex linguistic systems

- General learnability framework
- Previous learning successes: biased learners
- Unbiased probabilistic learning
- Where the problem lies: tricky data

Where next? Implications & Extensions



Where we are now

Modeling: aimed at understanding how children learn language, generating child behavior by using **psychologically plausible** methods



Learning complex systems: difficult.
Correct grammar is not the optimal grammar for child's input data without some kind of additional bias.

Where we are now

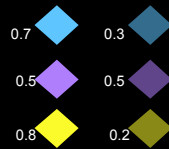
Modeling: aimed at understanding how children learn language, generating child behavior by using **psychologically plausible** methods



Learning complex systems: difficult.
Correct grammar is not the optimal grammar for child's input data without some kind of additional bias.

Bias on hypothesis space:
linguistic parameters already known, some values already known

Bias on data (Pearl 2008):
interpretive bias to use highly informative data



Where we can go

(1) Interpretive bias:

How successful on other difficult learning cases (noisy data sets, other complex systems)?

Are there other methods of implementing interpretative biases that lead to successful learning (productive data: Yang 2005)?

How necessary is an interpretive bias? Are there other biases that might cause the correct grammar to be the optimal grammar for the English data?



+ biases?

Where we can go

(1) Interpretive bias:

How successful on other difficult learning cases (noisy data sets, other complex systems)?

Are there other methods of implementing interpretative biases that lead to successful learning (productive data: Yang 2005)?

How necessary is an interpretive bias? Are there other biases that might cause the correct grammar to be the optimal grammar for the English data?



+ biases?

(2) Hypothesis space bias:

Will other hypothesis space instantiations allow the correct grammar to be the optimal grammar (constraints (Tesar & Smolensky 2000))? What learning mechanisms make the correct grammar learnable in these hypothesis spaces?

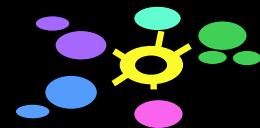
Is it possible to converge on the correct grammar given a less well-defined hypothesis space a priori (e.g. only knowing that units larger than syllables are required)?



+ other/fewer biases?

The big idea

Complex linguistic systems may well require something beyond probabilistic methods in order to be learned as well as children learn them given the data children are given.



What this likely is: learner biases in hypothesis space and data intake (how to deploy probabilistic learning)

What we can do with computational modeling:

(a) empirically test learning strategies that would be difficult to investigate with standard techniques

(b) generate experimentally testable predictions about learning (Pearl 2008: learning trajectory)



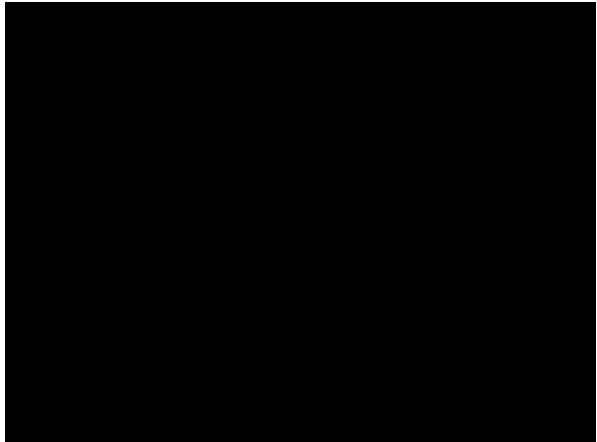
Thank You

Amy Weinberg
Bill Idsardi
Bill Sakas

Jeff Lidz
Charles Yang
Janet Fodor

The audiences at

University of California, San Diego Linguistics Department
UC Irvine Machine Learning Group
University of California, Los Angeles Linguistics Department
University of Southern California Linguistics Department



A Brief Tour of Parametric Metrical Phonology

Are syllables differentiated?

No: system is quantity-insensitive (QI)

S	S	S
CVV	CV	CCVC
lu	di	crous

A Brief Tour of Parametric Metrical Phonology

Are syllables differentiated?

No: system is quantity-insensitive (QI)

S	S	S
CVV	CV	CCVC
lu	di	crous

Yes: system is quantity-sensitive (QS)

Only allowed method: differ by rime weight

S	S	S
CVV	CV	CCVC
lu	di	crous

krəs
crous

Syllable

onset rime

kr /

 / \

 nucleus coda

 a s

A Brief Tour of Parametric Metrical Phonology

Are syllables differentiated?

No: system is quantity-insensitive (QI)

CVV CV CCVC
lu di crous

Yes: system is quantity-sensitive (QS)

Only allowed method: differ by rime weight
Only allowed number of divisions: 2
Heavy vs. Light

VV always Heavy
V always Light

Option 1: VC Heavy (QS-VC-H)

H L H
CVV CV CCVC
lu di crous

Option 2: VC Light (QS-VC-L)

H L L
CVV CV CCVC
lu di crous

narrowing of hypothesis space

A Brief Tour of Parametric Metrical Phonology

Are all syllables included in metrical feet?

Yes: system has no extrametricality (Em-None)

(L ... H)
VC VC VV
af ter noon

A Brief Tour of Parametric Metrical Phonology

Are all syllables included in metrical feet?

Yes: system has no extrametricality (Em-None)

(L ... H)
VC VC VV
af ter noon

No: system has extrametricality (Em-Some)

Only allowed # of exclusions: 1
Only allowed exclusions:
Leftmost or Rightmost syllable

narrowing of hypothesis space

A Brief Tour of Parametric Metrical Phonology

Are all syllables included in metrical feet?

Yes: system has no extrametricality (Em-None)

(L ... H)
VC VC VV
af ter noon

No: system has extrametricality (Em-Some)

Only allowed # of exclusions: 1
Only allowed exclusions:
Leftmost or Rightmost syllable

Leftmost syllable excluded: Em-Left
(...)
L H L
V VC V
a gen da

Rightmost syllable excluded: Em-Right
(...)
H L H
VV V VC
lu di crous

narrowing of hypothesis space

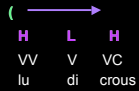
A Brief Tour of Parametric Metrical Phonology

What direction are metrical feet constructed?

Two logical options

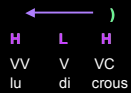
From the left:

Metrical feet are constructed from the left edge of the word (Ft Dir Left)



From the right:

Metrical feet are constructed from the right edge of the word (Ft Dir Right)



A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?

Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).

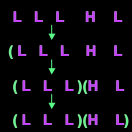
↑
narrowing of hypothesis space

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?

Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).

Ft Dir Left →



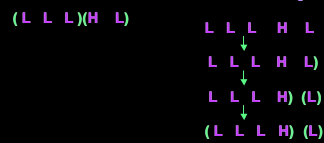
A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?


Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).

Ft Dir Left →

← Ft Dir Right



A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size? 


Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).

Ft Dir Left → (L L L)(H L) ← Ft Dir Right (L L L H)(L)

Ft Dir Left/Right
(L L L L L)
↓
(L L L L L)

(S S S S S)
↓
(S S S S S)

A Brief Tour of Parametric Metrical Phonology


Are metrical feet unrestricted in size?  (L L L)(H L)

Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**). (L L L H)(L)
(L L L L L)
(S S S S S)

No: Metrical feet are restricted (**Bounded**).

The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?  (L L L)(H L)

Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**). (L L L H)(L)
(L L L L L)
(S S S S S)

No: Metrical feet are restricted (**Bounded**).


The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space

Ft Dir Left → 2 units per foot (**Bounded-2**) 3 units per foot (**Bounded-3**)

x x x x
↓
(x x)(x x)
↓
(x x)(x x)

x x x x
↓
(x x x)(x)
↓
(x x x)(x)

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?  (L L L)(H L)

Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**). (L L L H)(L)
(L L L L L)
(S S S S S)

No: Metrical feet are restricted (**Bounded**).

The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space
The counting units are restricted to 2 options: syllables or moras.

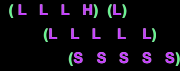
(x x)(x x) B-2
(x x x)(x) B-3

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?

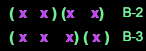


Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).



No: Metrical feet are restricted (**Bounded**).

The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space
The counting units are restricted to 2 options: syllables or moras.



Ft Dir Left Bounded-2

→ (H L)(L H)

← Count by syllables (Bounded-Syllabic)

(L L)(L H)

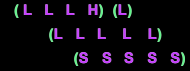
(S S)(S S)

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?

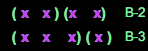


Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).



No: Metrical feet are restricted (**Bounded**).

The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space
The counting units are restricted to 2 options: syllables or moras.



Count by syllables (Bounded-Syllabic)

→ (H L)(L H)

Ft Dir Left Bounded-2

→ Count by moras (Bounded-Moraic)

xx x x xx

H L L H

↓

(H)(L L)(H)

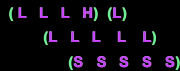
Moras (unit of weight):
H = 2 moras xx
L = 1 mora x

A Brief Tour of Parametric Metrical Phonology

Are metrical feet unrestricted in size?

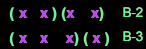


Yes: Metrical feet are unrestricted, delimited only by Heavy syllables if there are any (**Unbounded**).



No: Metrical feet are restricted (**Bounded**).

The size is restricted to 2 options: 2 or 3. ← narrowing of hypothesis space
The counting units are restricted to 2 options: syllables or moras.



Count by syllables (Bounded-Syllabic)

(H L)(L H)

Ft Dir Left Bounded-2

→ Count by moras (Bounded-Moraic)

(H)(L L)(H)

← compare

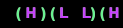
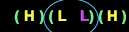
A Brief Tour of Parametric Metrical Phonology

Within a metrical foot, which syllable is stressed?

Two options, hypothesis space restriction

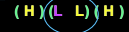
Leftmost:

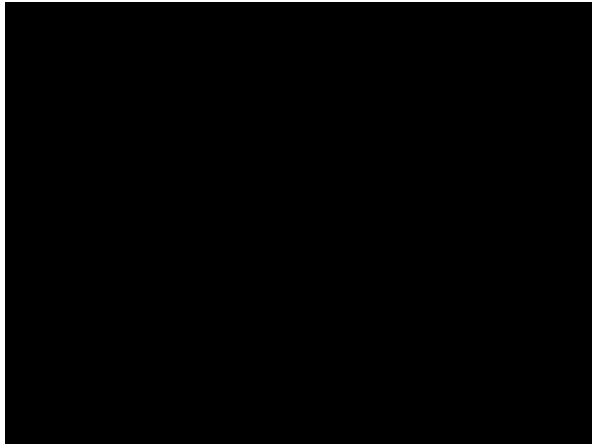
Stress the leftmost syllable (Ft Hd Left)



Rightmost:

Stress the rightmost syllable (Ft Hd Right)





Choosing among grammars

2ⁿ options

Human learning seems to be gradual and somewhat robust to noise - need some **probabilistic learning component**

Since grammars are parameterized, child can make use of this information to constrain hypothesis space. Learn over parameters, not entire parameter value sets.

probabilistic learning over parameter values

2n options

A caveat about learning parameters separately

or ? Parameters are system components that combine together to generate output.

or ? Choice of one parameter may influence choice of subsequent parameters.

or ?

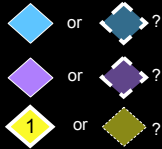
A caveat about learning parameters separately

or ? Parameters are system components that combine together to generate output.

or ? Choice of one parameter may influence choice of subsequent parameters.

or ?

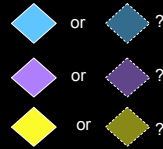
A caveat about learning parameters separately



Parameters are system components that combine together to generate output.

Choice of one parameter may influence choice of subsequent parameters.

A caveat about learning parameters separately



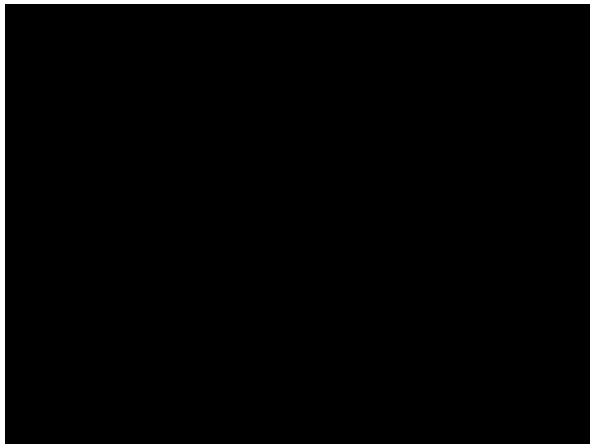
Parameters are system components that combine together to generate output.

Choice of one parameter may influence choice of subsequent parameters.



Dresher 1999

Point: The order in which parameters are set may determine if they are set correctly from the data.



Practical matters: Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Cues (Dresher, 1999; Lightfoot, 1999)



Parsing (Fodor, 1998; Sakas & Fodor, 2001)



Both operate over a single data point at a time:
compatible with incremental learning

Probabilistic learning from unambiguous data

(Pearl 2008)

Each parameter has 2 values.



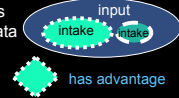
Probabilistic learning from unambiguous data

(Pearl 2008)

Each parameter has 2 values.



Advantage in data: How much more unambiguous data there is for one value over the other in the data distribution.



Assumption (Yang 2002):
The value with the greater advantage will be the one a probabilistic learner will converge on over time.

Allows us to be fairly agnostic about the exact nature of the probabilistic learning, provided it has this behavior.

Probabilistic learning from unambiguous data

(Pearl 2008)



The order in which parameters are set may determine if they are set correctly from the data.

Dresher 1999

Probabilistic learning from unambiguous data

(Pearl 2008)



The order in which parameters are set may determine if they are set correctly from the data.

Dresher 1999

Success guaranteed as long as parameter-setting order constraints are followed.

Cues

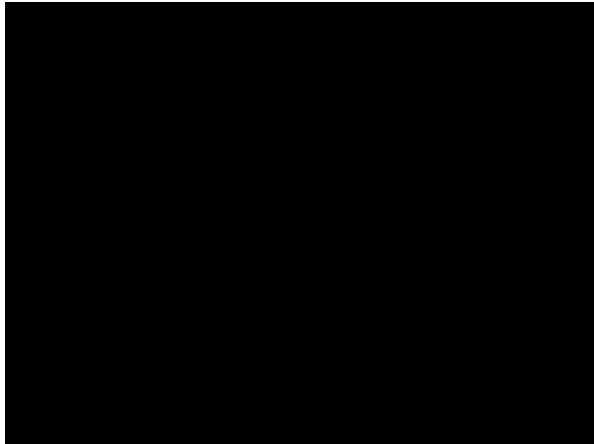
- (a) QS-VC-Heavy
before Em-Right
- (b) Em-Right
before Bounded-Syl
- (c) Bounded-2
before Bounded-Syl

The rest of the parameters are freely ordered w.r.t. each other.

Parsing

- Group 1:
QS, Ft Hd Left, Bounded
- Group 2:
Ft Dir Right, QS-VC-Heavy
- Group 3:
Em-Some, Em-Right, Bounded-2, Bounded-Syl

The parameters are freely ordered w.r.t. each other within each group.



Practical matters:
Feasibility of unambiguous data


Existence? Depends on data set (empirically determined).


Practical matters:
Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Cues (Dresher, 1999; Lightfoot, 1999) 

Parsing (Fodor, 1998; Sakas & Fodor, 2001) 

Both operate over a single data point at a time:
compatible with incremental learning

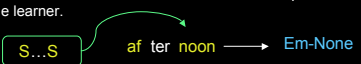
Practical matters:
Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Cues (Dresher 1999; Lightfoot 1999): heuristic pattern-matching to observable form of the data. Cues are available for each parameter value, known already by the learner.



Practical matters: Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Cues (Dresher 1999; Lightfoot 1999): **heuristic pattern-matching to observable form of the data.** Cues are available for each parameter value, known already by the learner.

QS: 2 syllable word with 2 stresses

VV VV

Em-Right: Rightmost syllable is Heavy and unstressed

... H

Unb: 3+ unstressed S/L syllables in a row

... S S S ...

... L L L L

Ft Hd Left: Leftmost foot has stress on leftmost syllable

S S S ...

H L L ...

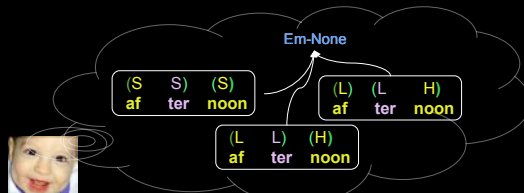
Practical matters: Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Parsing (Fodor 1998; Sakas & Fodor 2001): **extract necessary parameter values from all successful parses of data point** (strongest form of parsing)



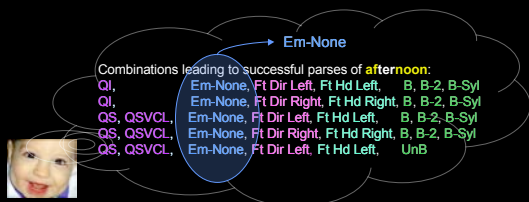
Practical matters: Feasibility of unambiguous data

Existence? Depends on data set (empirically determined).

Identification?

Identifying unambiguous data:

Parsing (Fodor 1998; Sakas & Fodor 2001): **extract necessary parameter values from all successful parses of data point** (strongest form of parsing)



Probabilistic learning for English: Modifications

Probabilistic generation and testing of parameter values (Yang 2002)

Learner bias: metrical phonology relies in part on knowledge of rhythmical properties of the language

Human infants may already have knowledge of Ft Hd Left and QS.

Build this bias into a model: set probability of QS = Ft Hd Left = 1.0.
These will always be chosen during generation.

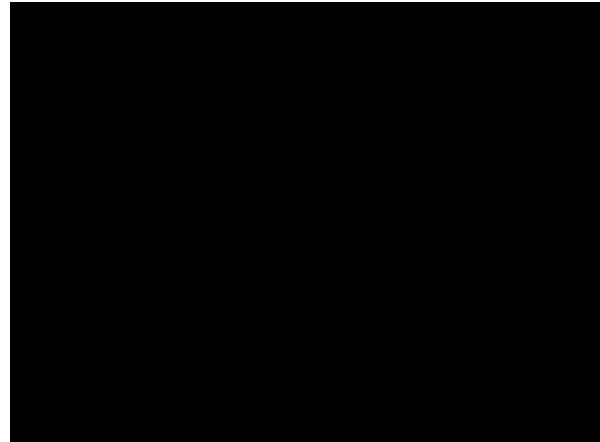
QS...QSVCL or QSVCH?

...
Ft Hd Left



QS, QSVCL, Em-None, Ft Dir Right,
Bounded, Bounded-2, Bounded-Syl, Ft Hd Left

Update parameter value probabilities + Batch Learning



Initial State of English Child-Directed Speech: Probability of Encountering Unambiguous Data

QS more probable

Em-None more probable

Quantity Sensitivity		Extrametricity	
QI: .00398	QS: 0.0205	None: 0.0294	Some: .0000259
Feet Directionality		Boundedness	
Left: 0.000	Right: 0.00000925	Unbounded: 0.00000370	Bounded: 0.00435
Feet Headedness			
Left: 0.00148	Right: 0.000		

Moving Targets & Unambiguous Data: What Happens After Parameter-Setting

Em-None more probable

Quantity Sensitivity		Extrametricity	
QI: .00398	QS: 0.0205	None: 0.0294	Some: .0000259
Feet Directionality		Boundedness	
Left: 0.000	Right: 0.00000925	Unbounded: 0.00000370	Bounded: 0.00435
Feet Headedness			
Left: 0.00148	Right: 0.000		

Moving Targets & Unambiguous Data: What Happens After Parameter-Setting

Em-Some more probable

QS		Extrametricality	
		None: 0.0240	Some: .0485
Feet Directionality		Boundedness	
Left: 0.000	Right: 0.00000555	Unbounded: 0.00000370	Bounded: 0.00125
Feet Headedness			
Left: 0.000588	Right: 0.0000204		

Getting to English

The child must set all the parameter values in order to converge on a language system.

Current knowledge of the system (parameters set) influences the perception of unambiguous data (subsequent parameters set).

QS
?



Dresher 1999

The order in which parameters are set may determine if they are set correctly from the data.

Will any parameter-setting orders lead the learner to English?

Feasibility & Sufficiency of the Unambiguous Data Filter

Either method of identifying unambiguous data (cues or parsing) is **successful**. Given the **non-trivial parametric system** (9 interactive parameters) and the non-trivial data set (English is full of exceptions), this is no small feat.



Clark 1994

Existence?

"It is unlikely that any example ... would show the effect of only a single parameter value; rather, each example is the result of the interaction of several different principles and parameters"

Feasibility & Sufficiency of the Unambiguous Data Filter

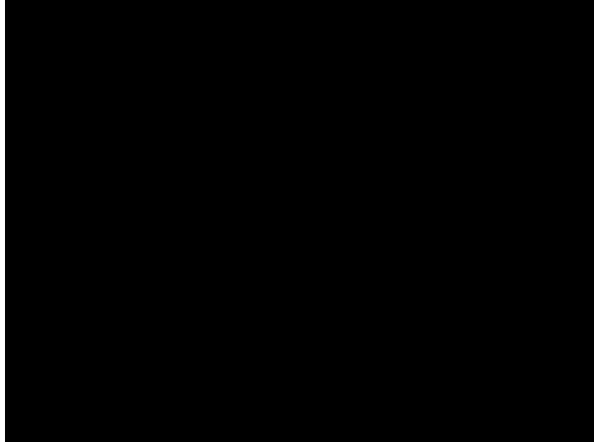
Either method of identifying unambiguous data (cues or parsing) is **successful**. Given the **non-trivial parametric system** (9 interactive parameters) and the non-trivial data set (English is full of exceptions), this is no small feat.

✓ Existence

✓ Identification

(1) **Unambiguous data** exist and can be identified in sufficient relative quantities to learn a **complex parametric system**.

(2) The **selective learning strategy** is robust across a realistic (highly ambiguous, exception-filled) data set. It's feasible to identify such data, and the strategy yields sufficient learning behavior.



Where we can go: Links to the Experimental Side

Cues

- (a) QS-VC-Heavy
before Em-Right
- (b) Em-Right
before Bounded-Syl
- (c) Bounded-2
before Bounded-Syl

Parsing

- Group 1:
QS, Ft Hd Left, Bounded
- Group 2:
Ft Dir Right, QS-VC-Heavy
- Group 3:
Em-Some, Em-Right, Bounded-2,
Bounded-Syl

Are predicted parameter setting orders observed in real-time learning?

E.g. whether cues or parsing is used, Quantity Sensitivity (QS, QSVCH) is predicted to be set before Extrametricality (Em-Some, Em-Right).

And in fact, there is evidence that quantity sensitivity may be known quite early (Turk, Jusczyk, & Gerken, 1995)