

# Hierarchical Bayesian Models for Modeling Cognitive Processes

Çağrı Çöltekin

Center for Language and Cognition  
University of Groningen  
c.coltekin@rug.nl

March 11, 2009

# Overview

## Bayesian Inference

- Introduction (with comparison to frequentist statistics)

- A Simple Example

## Bayesian Inference and Learning

- MLE and MAP

- Bayesian Learning

## Hierarchical Bayesian Models

- A simple HBM example

## Using HBMs for Learning Grammar

- Categorical Grammars

- A Bayesian CG learner

- A hierarchical extension

# Two views on statistics

## ▶ **Frequentist view**

- ▶ Probabilities are long-run frequencies.
- ▶ We can talk about probabilities only for well defined experiments with random outcomes.
- ▶ Emphasis is on *objectivity*: data must speak for itself.

## ▶ **Bayesian view**

- ▶ Probabilities are degrees of belief.
- ▶ Probabilities can be assigned to any statement.
- ▶ Inference is *subjective* (methods for somewhat objective inference exists)

# Statistical Inference

- ▶ We collect a sample  $\mathbf{X}$  from a population.
- ▶ We know (or assume) that  $\mathbf{X}$  is according to a known distribution that can be parametrized by  $\theta$

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta)p(\theta)}{p(\mathbf{X})}$$

$p(\theta|\mathbf{X})$  : *posterior*

$p(\mathbf{X}|\theta)$  : *likelihood* ( $\mathcal{L}(\theta)$ )

$p(\theta)$  : *prior*

$p(\mathbf{X})$  : Marginal probability of data ( $\int p(\mathbf{X}|\theta)p(\theta)d\theta$ )

# Statistical Inference: the two approaches

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta)p(\theta)}{p(\mathbf{X})}$$
$$\propto p(\mathbf{X}|\theta)p(\theta)$$

- ▶ Frequentist approach:
  - ▶  $\theta$  is a fixed but unknown.
  - ▶ Inference is done via *Maximum Likelihood Estimate* (MLE).  
Prior information is never used.
  - ▶ We need to assess the reliability of the estimate by significance tests (p values, confidence intervals).
- ▶ Bayesian approach:
  - ▶  $\theta$  is treated just like any other random variable.
  - ▶ Posterior contains all the necessary ingredients for the inference.

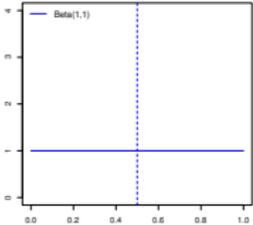
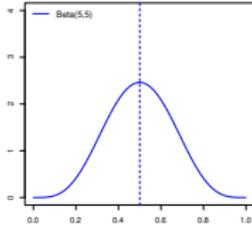
## A simple example

We toss a coin 10 times, the outcome is *HHTHHHHTTH*. Let  $\theta$  represent the chance that coin comes up 'H'.

Frequentist	Bayesian 1	Bayesian 2

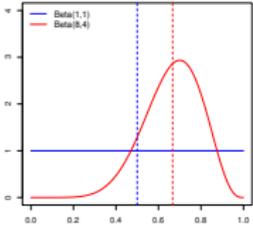
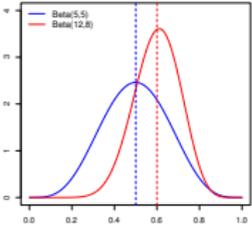
## A simple example

We toss a coin 10 times, the outcome is *HHTHHHHTTH*. Let  $\theta$  represent the chance that coin comes up 'H'.

Frequentist	Bayesian 1	Bayesian 2
<i>No prior</i>	<p>Prior:</p> 	<p>Prior:</p> 

## A simple example

We toss a coin 10 times, the outcome is *HHTHHHHTTH*. Let  $\theta$  represent the chance that coin comes up 'H'.

Frequentist	Bayesian 1	Bayesian 2
$\hat{\theta} = \operatorname{argmax}_{\theta} \mathcal{L}(\theta)$ $\hat{\theta} = 0.7$ $p = 0.2059$	<p>Posterior:</p>  <p>The plot shows two probability density functions on the interval [0, 1]. A blue curve represents the uniform distribution Beta(1,1), which is a flat horizontal line at height 1. A red curve represents the Beta(8,4) distribution, which is a smooth curve peaking at approximately 0.67. Vertical dashed lines are drawn at 0.5 (blue) and 0.67 (red).</p>	<p>Posterior:</p>  <p>The plot shows two probability density functions on the interval [0, 1]. A blue curve represents the Beta(5,5) distribution, which is a smooth curve peaking at 0.5. A red curve represents the Beta(12,8) distribution, which is a smooth curve peaking at approximately 0.6. Vertical dashed lines are drawn at 0.5 (blue) and 0.6 (red).</p>

# Bayesian Inference: summary

$$p(\theta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta)$$

## Pros

- ▶ Allows complete inference: posterior distribution contains all the information needed.
- ▶ Interpretation of the results are straightforward.

## Cons

- ▶ Computationally expensive: calculation of posteriors are not always easy.
- ▶ Choice of priors: it is sometimes difficult to justify the use of (subjective) priors.

# Statistical Inference and Learning

- ▶ Statistical inference aims to draw conclusions about an underlying population using a sample drawn from this population.
- ▶ Learning can be viewed as making generalizations about the target concept by looking at a sample of data consistent with this concept.
- ▶ Statistical methods are most common learning methods in machine learning.
- ▶ More and more psychological phenomena is explained by sensitivity to statistical information in the environment.
- ▶ Language acquisition is no exception: children are known to exploit statistical regularities in the input in language learning.

# (non-Bayesian) Statistical Learning

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta)p(\theta)}{p(\mathbf{X})}$$
$$\propto p(\mathbf{X}|\theta)p(\theta)$$

We want to learn the parameter  $\theta$ .

- ▶ Maximum Likelihood Estimate(MLE):

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathbf{X}|\theta)$$

- ▶ Maximum a posteriori (MAP) estimate:

$$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathbf{X}|\theta)p(\theta)$$

# Bayesian Learning

$$p(\theta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta)$$

- ▶ No point estimates.
- ▶ No maximization.
- ▶ A Bayesian learner learns the posterior distribution,  $p(\theta|\mathbf{X})$ .
- ▶ If needed, point estimates can be made using,

$$E[\theta] = \int \theta p(\theta|\mathbf{X}) d\theta$$

Note the difference from MAP estimate.

# Bayesian Learning: how?

$$p(\theta|\mathbf{X}) = \frac{p(\mathbf{X}|\theta)p(\theta)}{\int p(\mathbf{X}|\theta)p(\theta)d\theta}$$

- ▶ Given the probability density (or distribution) functions for prior and likelihood, we simply multiply, and normalize.
- ▶ Note that likelihood and prior does not have to be *proper*. Multiplication by a constant does not change the results.

## Problem 1: computation

The computation involved is not always easy to carry out. Approximate methods, such as MCMC, are frequently used.

## Problem 2: priors

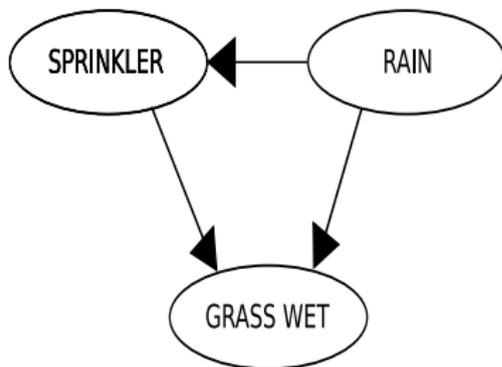
We need to choose a prior distribution.

# Choice of priors

- ▶ Subjective priors: based on previous experience.
- ▶ Non-informative priors: try to be as objective as possible.
- ▶ Conjugate priors: for computational efficiency.
- ▶ Empirical Bayes: priors from data.
- ▶ Hierarchical priors: combining information from different variables.

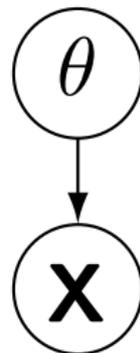
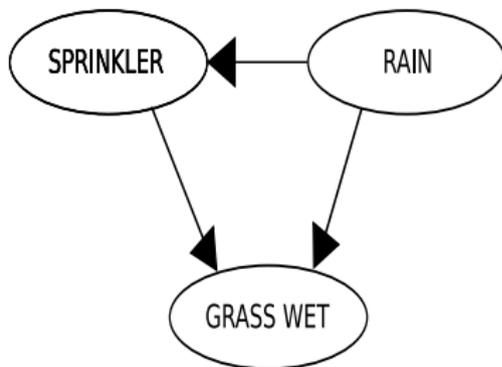
## Digression: Graphical models (or 'Bayesian networks')

- ▶ Often multiple random variables interact.
- ▶ Bayesian networks are a convenient way to visualize the dependency relations between variables.



## Digression: Graphical models (or 'Bayesian networks')

- ▶ Often multiple random variables interact.
- ▶ Bayesian networks are a convenient way to visualize the dependency relations between variables.



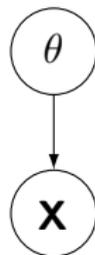
# Hierarchical Bayesian Models

In our coin toss example we choose a  $Beta(\alpha, \beta)$  prior with fixed  $\alpha$  and  $\beta$ .

$$p(\theta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta)$$

where,

$$\mathbf{X} \sim \text{Binomial}(\theta), \theta \sim \text{Beta}(\alpha, \beta)$$



# Hierarchical Bayesian Models

In our coin toss example we choose a  $Beta(\alpha, \beta)$  prior with fixed  $\alpha$  and  $\beta$ .

$$p(\theta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta)$$

where,

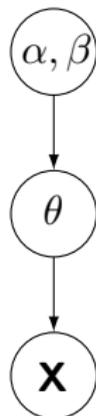
$$\mathbf{X} \sim \text{Binomial}(\theta), \theta \sim \text{Beta}(\alpha, \beta)$$

We can further extend this if choice of  $\alpha$  and  $\beta$  can be guided by additional information.

$$p(\theta, \alpha, \beta|\mathbf{X}) \propto p(\mathbf{X}|\theta)p(\theta|\alpha, \beta)p(\alpha, \beta)$$

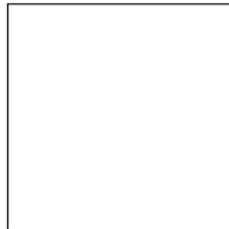
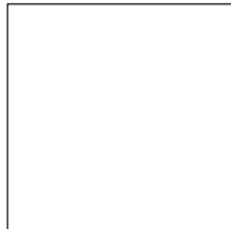
where, e.g.,

$$\mathbf{X} \sim \text{Binomial}(\theta), \theta \sim \text{Beta}(\alpha, \beta), \alpha \sim N(\mu_\alpha, \sigma), \beta \sim N(\mu_\beta, \sigma)$$



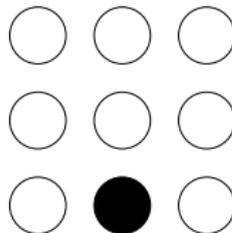
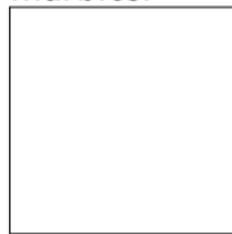
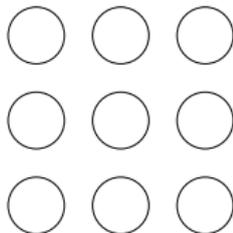
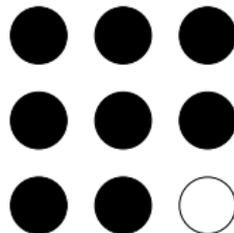
# A simple example: marbles in boxes

Boxes contain either white or black marbles:



## A simple example: marbles in boxes

Boxes contain either white or black marbles:



\* This example is adopted from a talk by J. Tenenbaum

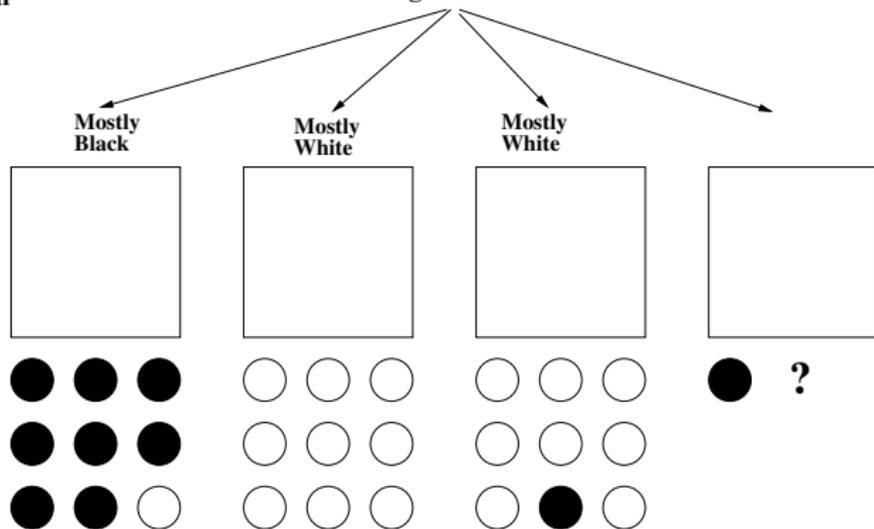
# A simple example: marbles in boxes

Boxes contain either white or black marbles:

Level 2: Boxes in general

Color varies among boxes not much within boxes

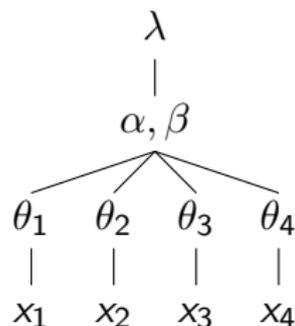
Level 1: Box properties



\* This example is adopted from a talk by J. Tenenbaum

# HBM for marbles in boxes

- ▶  $x_i$  are the data
- ▶  $\theta_i$  models box proportions,  
 $x_i \sim \text{Binomial}(\theta)$
- ▶  $\alpha, \beta$  models the boxes in general,  
 $\theta_i \sim \text{Beta}(\alpha, \beta)$
- ▶  $\lambda$  models prior expectations in boxes in general



$$p(\lambda, \alpha, \beta, \theta \mid \mathbf{X}) \propto p(\mathbf{X} \mid \theta)p(\theta \mid \lambda, \alpha, \beta)p(\alpha, \beta \mid \lambda)p(\lambda)$$

## Summary: HBMs for Modeling Human Learning

- ▶ Bayesian Statistics provides complete inference: posterior distribution contains all we need.
- ▶ Bayesian Learning is incremental: posterior can be used as prior for the next step.
- ▶ Hierarchical models allow a way to include information from different sources as prior knowledge.

# Overview

## Bayesian Inference

Introduction (with comparison to frequentist statistics)

A Simple Example

## Bayesian Inference and Learning

MLE and MAP

Bayesian Learning

## Hierarchical Bayesian Models

A simple HBM example

## Using HBMs for Learning Grammar

Categorical Grammars

A Bayesian CG learner

A hierarchical extension

# Categorial Grammars

- ▶ Categorial Grammar (CG) encodes all the language specific syntactic information in the lexicon.
- ▶ CG Lexicon contains *lexical items* of the form:

$$\phi := \sigma : \mu$$

where  $\phi$  is the *phonological form*,  $\sigma$  is the syntactic category, and  $\mu$  is the meaning of the lexical item.

- ▶ Syntactic categories in CG are,
  - ▶ either a basic category, such as **N**, **NP**, **S**
  - ▶ or, a complex category of the form **X\Y** or **X/Y**, where **X** and **Y** are any (basic or complex) CG categories. Informally:
    - X/Y** says: 'I need a **Y** to my *right* to become **X**'
    - X\Y** says: 'I need a **Y** to my *left* to become **X**'

# Categorial Grammars (contd.)

- ▶ More formally, CG has two rules:

Forward Application:  $X/Y \quad Y \Rightarrow \mathbf{X} \quad (>)$

Backward Application:  $Y \quad X \backslash Y \Rightarrow \mathbf{X} \quad (<)$

- ▶ An example derivation:

*Mary*      *likes*      *Peter*  
 $\overline{\text{NP}} \quad \overline{(\text{S} \backslash \text{NP}) / \text{NP}} \quad \overline{\text{NP}}$   
 $\xrightarrow{\text{S} \backslash \text{NP}}$   
 $\xleftarrow{\text{S}}$

# Categorial Grammars (contd.)

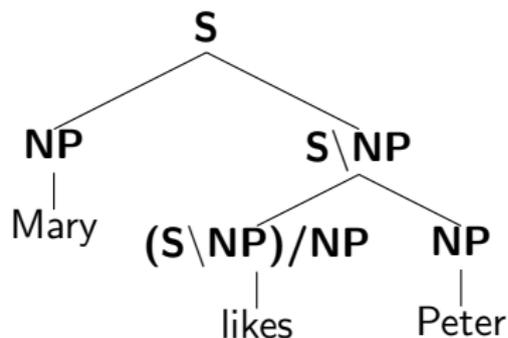
- ▶ More formally, CG has two rules:

Forward Application:  $X/Y \quad Y \Rightarrow \mathbf{X} \quad (>)$

Backward Application:  $Y \quad X \backslash Y \Rightarrow \mathbf{X} \quad (<)$

- ▶ An example derivation:

*Mary*      *likes*      *Peter*  
NP    (S \ NP) / NP    NP  
----->  
                  S \ NP  
-----<  
                  S



# A simple CG learner for learning word-grammars

- ▶ Input is a set of words.
- ▶ We want to assign a probability,  $\theta$ , to possible lexical items ( $\langle \phi, \sigma \rangle$  pairs).
- ▶ Note: probability is the 'system's belief' that the  $\langle \phi, \sigma \rangle$  pair at hand is a lexical item.
- ▶ Only 3 categories (known in advance):
  - ▶ **W** : free morpheme (word, or stem)
  - ▶ **W/W** :

# A simple CG learner for learning word-grammars

- ▶ Input is a set of words.
- ▶ We want to assign a probability,  $\theta$ , to possible lexical items ( $\langle \phi, \sigma \rangle$  pairs).
- ▶ Note: probability is the 'system's belief' that the  $\langle \phi, \sigma \rangle$  pair at hand is a lexical item.
- ▶ Only 3 categories (known in advance):
  - ▶ **W** : free morpheme (word, or stem)
  - ▶ **W/W** : prefix
  - ▶ **W\W** :

# A simple CG learner for learning word-grammars

- ▶ Input is a set of words.
- ▶ We want to assign a probability,  $\theta$ , to possible lexical items ( $\langle \phi, \sigma \rangle$  pairs).
- ▶ Note: probability is the 'system's belief' that the  $\langle \phi, \sigma \rangle$  pair at hand is a lexical item.
- ▶ Only 3 categories (known in advance):
  - ▶  $\mathbf{W}$  : free morpheme (word, or stem)
  - ▶  $\mathbf{W}/\mathbf{W}$  : prefix
  - ▶  $\mathbf{W}\backslash\mathbf{W}$  : suffix
- ▶ We adopt a *Beta/Binomial* model.
- ▶ We assume each input word provides evidence for the lexical hypothesis in question, If hypothesis used in the interpretation of the input.

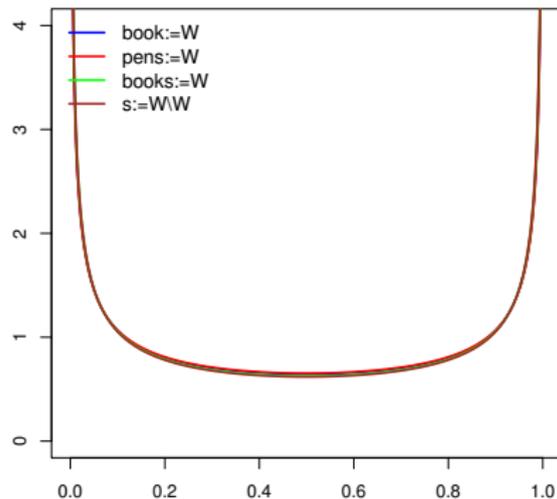
# The CG learner: A simple algorithm

- ▶ Input is unsegmented words (and the lexicon).
- ▶ Output is the lexicalized grammar with probability assignments.

For each input word  $w$ ,

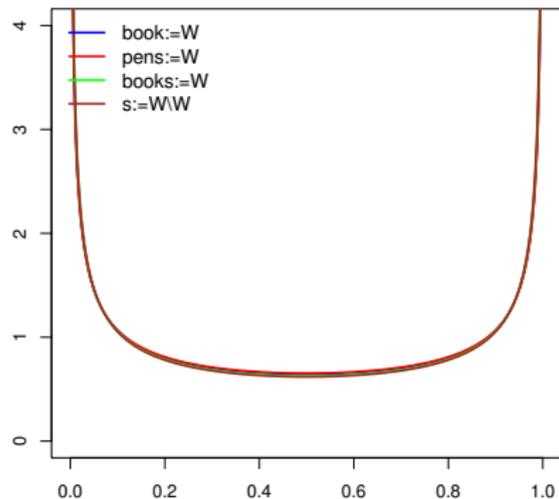
1. Try to segment the input using the current lexicon.
2. If there is no possible segmentation, assume that we have found evidence for a lexical item  $w := W$ .
3. If we can segment the input as  $w = \phi_1 \dots \phi_N$ , assume that we have observed evidence for each tuple  $\langle \phi_i, \sigma_j \rangle$  which yields a correct parse of  $w$ .
4. We update the parameters of the Beta distribution associated with the lexical hypotheses.

# An example



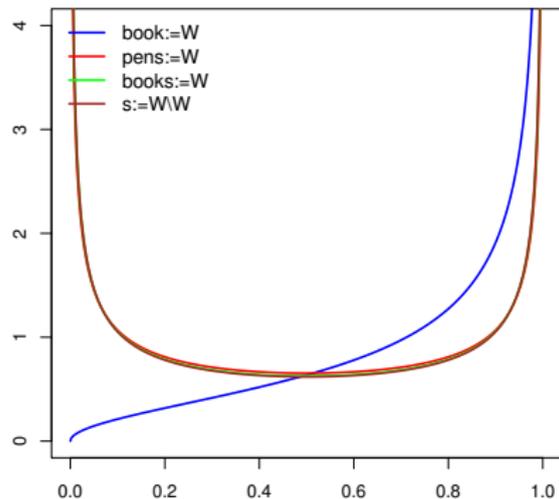
Lexicon      {  
Input        book  
Hypotheses  
Parses

# An example



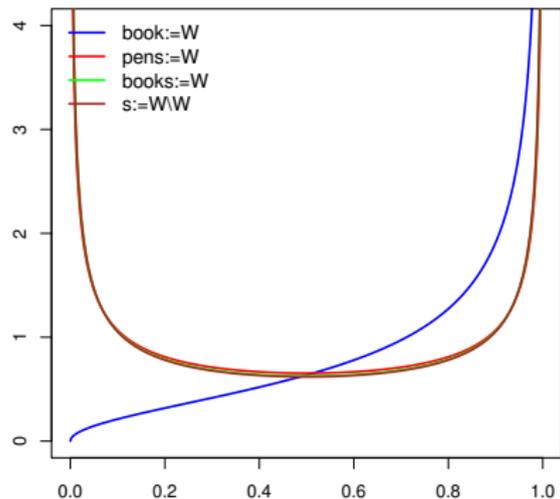
Lexicon	{ }
Input	book
Hypotheses	book:=W
Parses	

# An example



Lexicon	{book:=W}
Input	book
Hypotheses	book:=W
Parses	<b>book</b> <u>W</u>

# An example



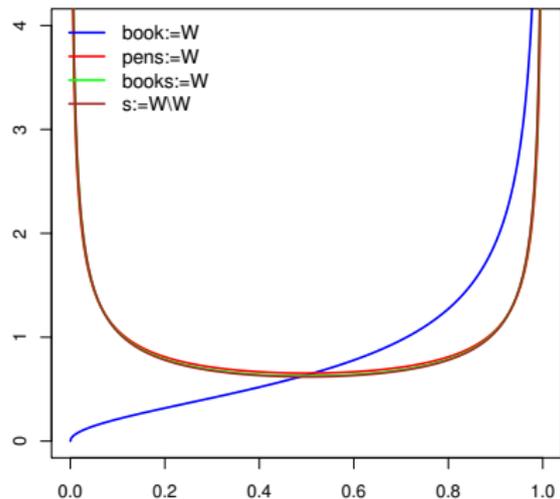
Lexicon      {**book:=W**}

Input        pens

Hypotheses

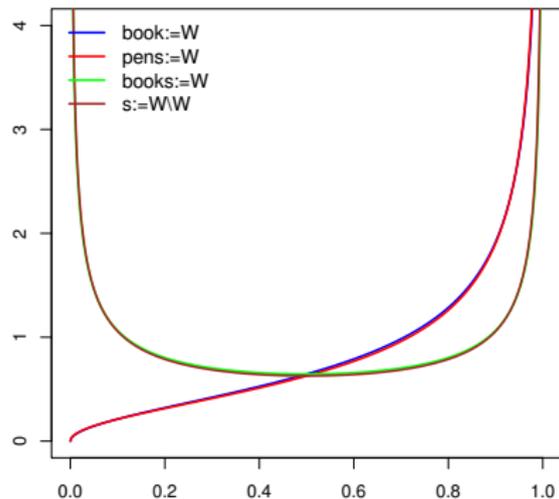
Parses

# An example



Lexicon	{book:=W}
Input	pens
Hypotheses	pens:=W
Parses	

# An example



Lexicon {book:=W,  
pens:=W}

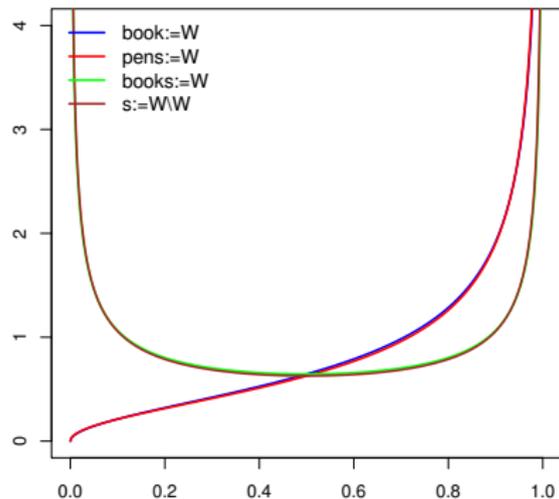
Input pens

Hypotheses pens:=W

Parses pens

W

# An example



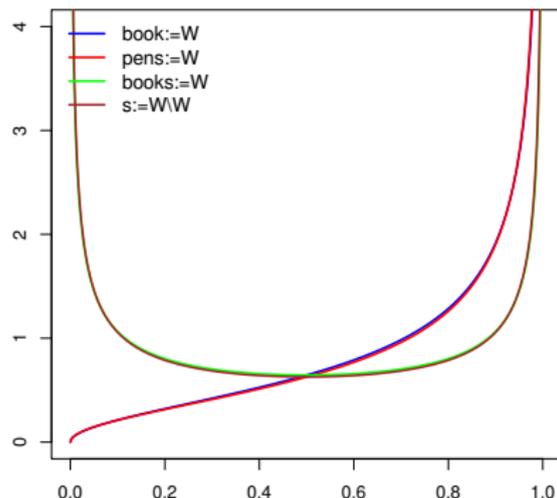
Lexicon {**book:=W**,  
**pens:=W**}

Input books

Hypotheses

Parses

# An example



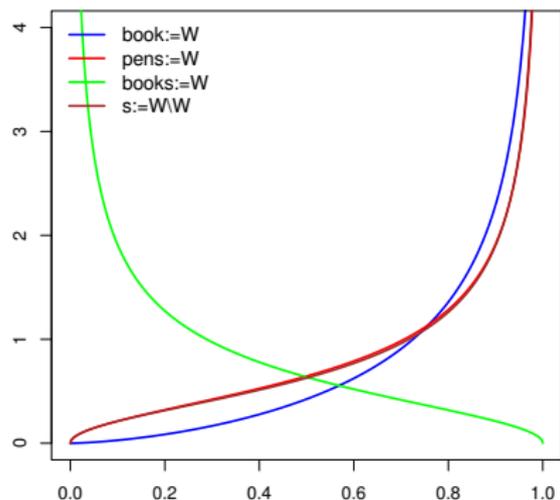
Lexicon {book:=W,  
pens:=W}

Input books

Hypotheses book:=W, s:=W,  
books:=W,  
book:=W/W,  
s:=W\W

Parses

# An example



Lexicon {book:=W,  
pens:=W, s:=W\W}

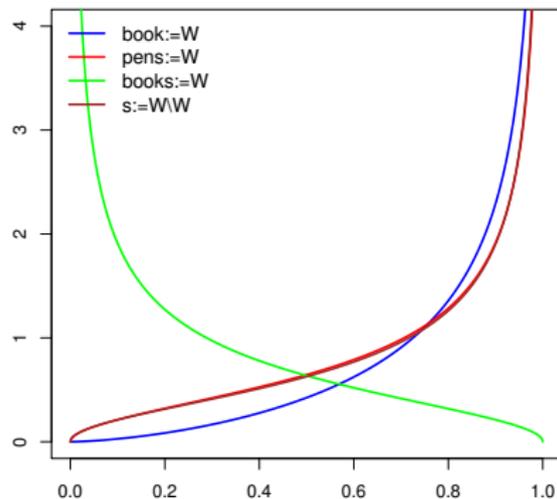
Input books

Hypotheses book:=W, s:=W,  
books:=W,  
book:=W/W,  
s:=W\W

Parses

books  
 $\overline{W}$   
**book** **s**  
 $\overline{W}$   $\overline{W \setminus W}$   
 $\overline{W}$  <  
 book s  
 $\overline{W/W}$   $\overline{W}$   
 $\overline{W}$  >

# An example



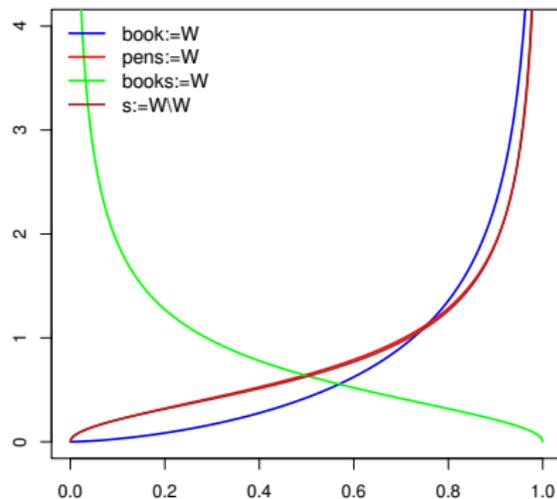
Lexicon

{**book**:=W,  
**pens**:=W, **s**:=W\W}

Input  
Hypotheses  
Parses

pens

# An example



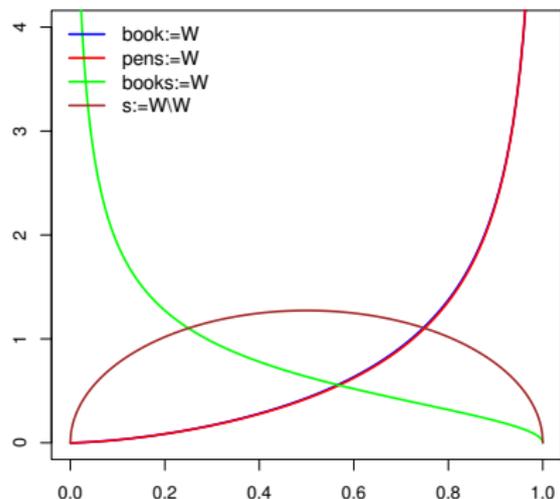
Lexicon {book:=W,  
pens:=W, s:=W\W}

Input pens

Hypotheses pen:=W, pens:=W,  
s:=W, pen:=W/W,  
s:=W\W

Parses

# An example



Lexicon  $\{\text{book}:=W, \text{pens}:=W, \text{s}:=W \setminus W\}$

Input pens

Hypotheses pen:=W, pens:=W, s:=W, pen:=W/W, s:=W \setminus W

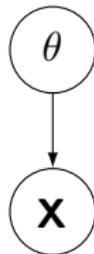
Parses

$\overline{\text{pens}}$       pen      s  
 $\overline{\text{W}}$        $\overline{\text{W}}$        $\overline{\text{W} \setminus \text{W}}$   
 $\overline{\text{W}}$  <

pen      s  
 $\overline{\text{W/W}}$   $\overline{\text{W}}$   
 $\overline{\text{W}}$  >

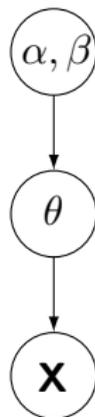
## A hierarchical extension

- ▶ Our current model assumes rather non-informative values for  $\alpha$  and  $\beta$ .
- ▶ We can extend this model to get more informative priors

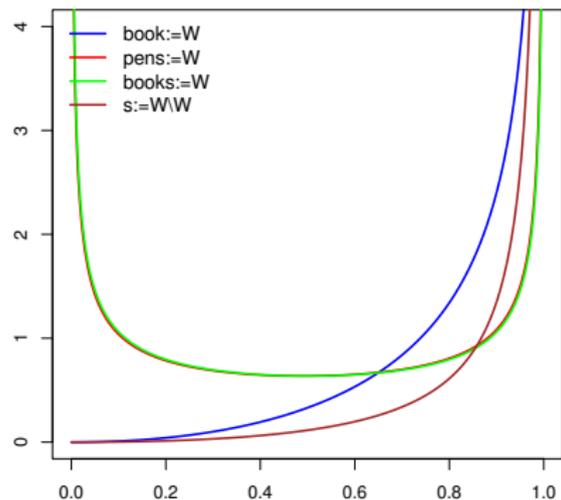


## A hierarchical extension

- ▶ Our current model assumes rather non-informative values for  $\alpha$  and  $\beta$ .
- ▶ We can extend this model to get more informative priors
- ▶ We treat  $\alpha$  and  $\beta$  as random variables.
- ▶ We make use of *context predictability* as another source providing a hierarchical informative prior for possible segments.

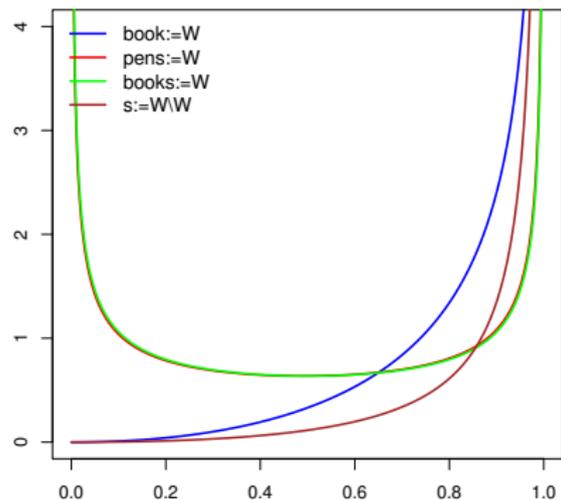


# Hierarchical extension: example



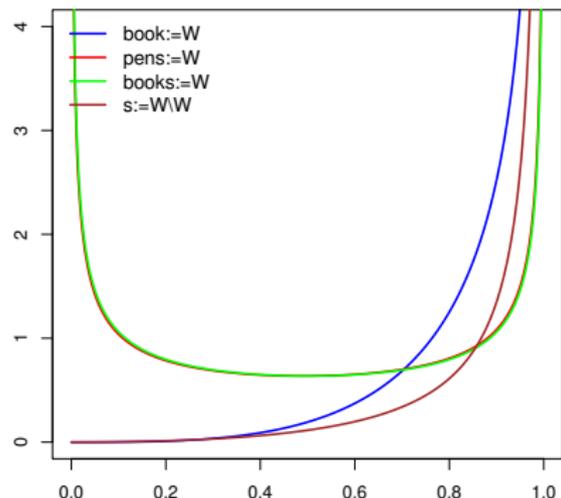
Lexicon	{
Input	book
Hypotheses	
Parses	

## Hierarchical extension: example



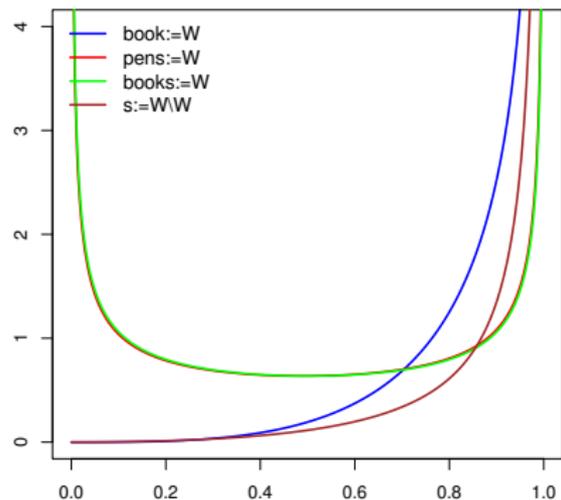
Lexicon	{}
Input	book
Hypotheses	book:=W
Parses	

# Hierarchical extension: example



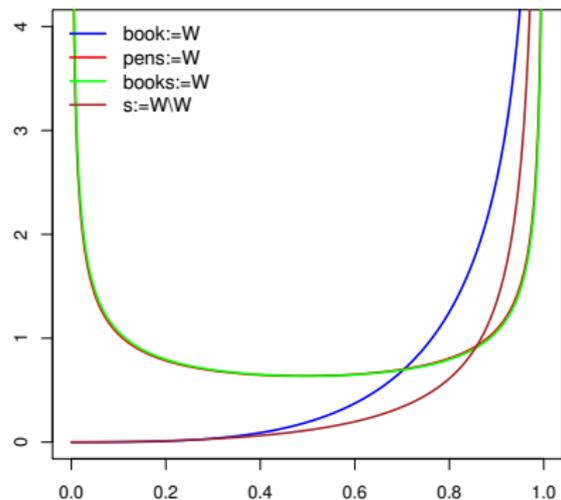
Lexicon	{book:=W}
Input	book
Hypotheses	book:=W
Parses	<b>book</b> <u>W</u>

# Hierarchical extension: example



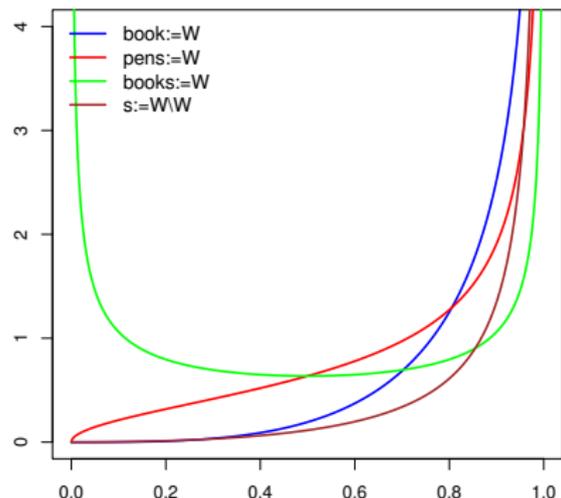
Lexicon	{book:=W}
Input	pens
Hypotheses	
Parses	

# Hierarchical extension: example



Lexicon	{book:=W}
Input	pens
Hypotheses	pens:=W
Parses	

# Hierarchical extension: example



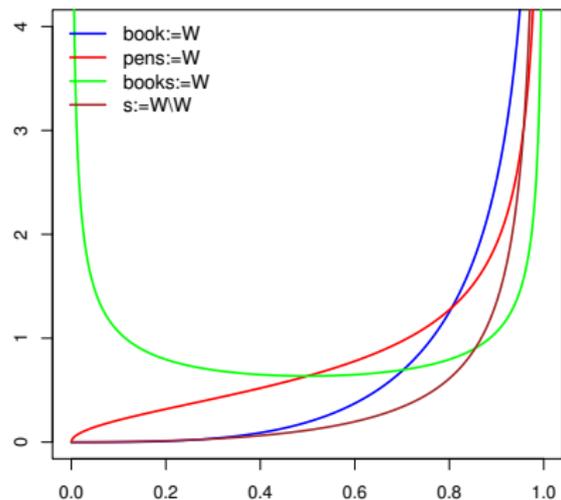
Lexicon {book:=W,  
pens:=W}

Input pens

Hypotheses pens:=W

Parses **pens**  
W

# Hierarchical extension: example



Lexicon      {book:=W,  
pens:=W}

---

Input        books

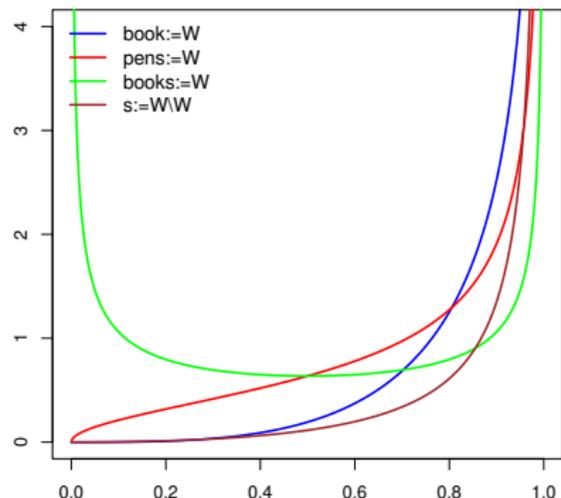
---

Hypotheses

---

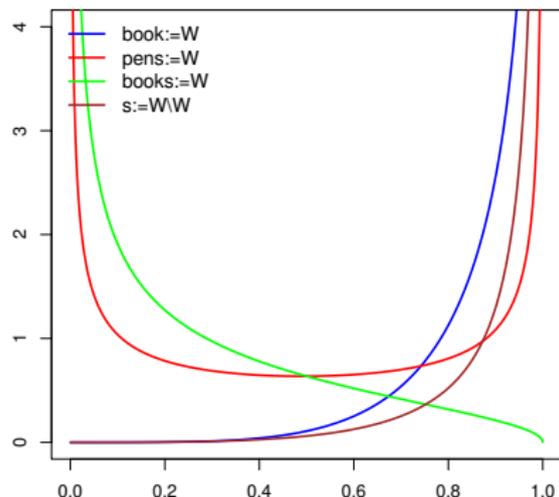
Parses

# Hierarchical extension: example



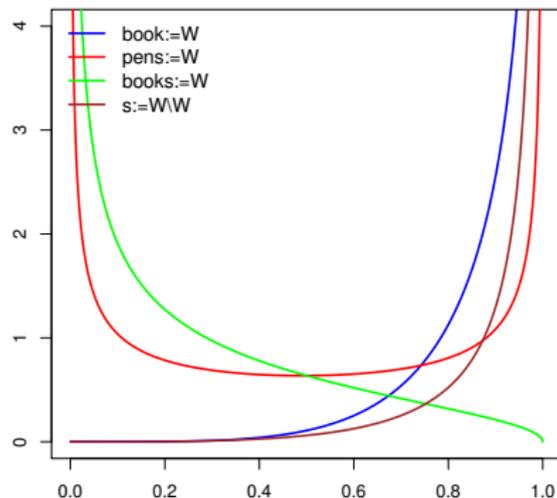
Lexicon	$\{\text{book}:=W,$ $\text{pens}:=W\}$
Input	books
Hypotheses	$\text{book}:=W,$ $\text{books}:=W,$ $s:=W,$ $\text{book}:=W/W,$ $s:=W \setminus W$
Parses	

# Hierarchical extension: example



Lexicon	{ <b>book:=W</b> , <b>pens:=W</b> , <b>s:=W\W</b> }
Input	books
Hypotheses	<b>book:=W</b> , <b>books:=W</b> , <b>s:=W</b> , <b>book:=W/W</b> , <b>s:=W\W</b>
Parses	books $\overline{W}$ <b>book</b> <b>s</b> $\overline{W}$ $\overline{W \setminus W}$ $\overline{W}$ < book   s $\overline{W/W}$ $\overline{W}$ $\overline{W}$ >

# Hierarchical extension: example



Lexicon      {book:=W,  
pens:=W, s:=W\W}

---

Input          pens

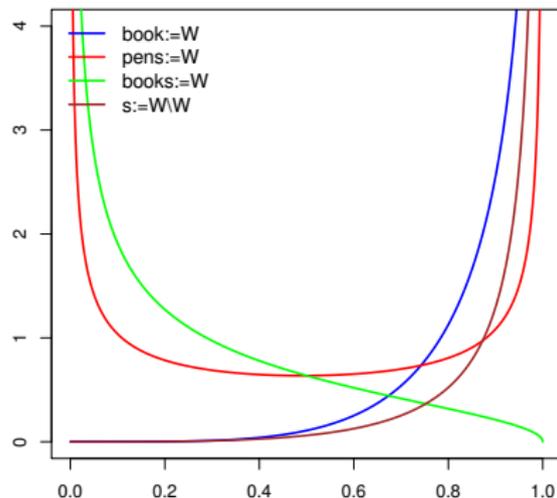
---

Hypotheses

---

Parses

# Hierarchical extension: example



Lexicon  $\{\text{book}:=W,$   
 $\text{pens}:=W, \text{s}:=W \setminus W\}$

---

Input pens

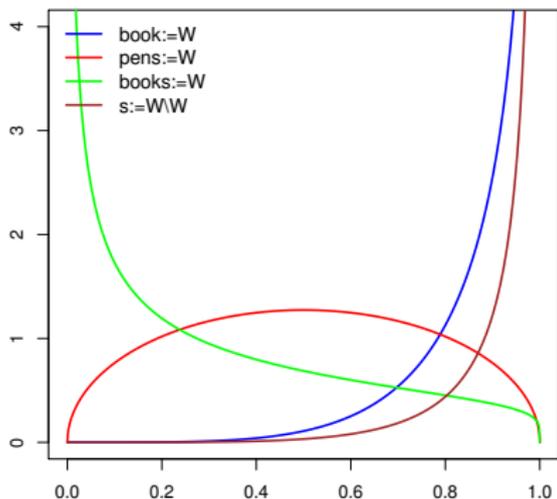
---

Hypotheses pen:=W, pens:=W,  
s:=W, pen:=W/W,  
s:=W \setminus W

---

Parses

# Hierarchical extension: example



Lexicon  $\{\text{book}:=W,$   
 $\text{pens}:=W, s:=W \setminus W,$   
 $\text{pen}:=W\}$

Input pens

Hypotheses  $\text{pen}:=W, \text{pens}:=W,$   
 $s:=W, \text{pen}:=W/W,$   
 $s:=W \setminus W$

Parses

$\text{pens}$	$\text{pen}$	$s$
$\overline{W}$	$\overline{W}$	$\overline{W \setminus W}$
	$\overline{W} <$	
$\text{pen}$	$s$	
$\overline{W/W}$	$\overline{W}$	
$\overline{W} >$		

# Summary

- ▶ Bayesian statistics provides a different approach to statistical inference and learning.
- ▶ Use of (subjective) priors is not always bad: Modeling cognitive processes is a good example.
- ▶ Hierarchical priors is a good way to combine information from different sources.