



Machine Learning for the Computational Humanities

David Bamman
Carnegie Mellon University

Oct 24, 2014

#mlch

Overview

- Classification
 - Probability
 - Independent (Logistic regression, Naive Bayes)
 - Structured (CRFs, HMMs)
- Clustering (hierarchical, K-means)
- Probabilistic graphical models (e.g., topic models)
- Representation learning

The big two

- Classification
 - Given a pre-defined set of categories, determine which category (or categories) apply to the text.
Example: spam vs. not spam.
- Clustering
 - Learn coherent groups according to some notion of similarity.

Classification

- Supervised classification learns a mapping from an input to an output from training data

Application	Input	Output
Spam filtering	email	spam, not spam
Authorship attribution	document	author
Sentiment analysis	text	positive, negative
Part of speech tagging	sentence	sequence of part of speech tags

Training data

Label	Input
Jane Austen	It is a truth universally acknowledged, that a single man in possession ...
Jane Austen	Emma Woodhouse, handsome, clever, and rich, with a comfortable home ...
Jane Austen	The family of Dashwood had long been settled in Sussex. Their estate...
Jane Austen	Sir Walter Elliot, of Kellynch Hall, in Somersetshire, was a man who, for ...
Herman Melville	Call me Ishmael. Some years ago--never mind how long precisely...
Herman Melville	I am a rather elderly man. The nature of my avocations for the last thirty ...
Mark Twain	You don't know about me without you have read a book by the name of...

What do you need?

Two steps to building and using a supervised classification model.

1. **Train** a model with data where you know the answers.
2. Use that model to **predict** data where you don't.

What do you need?

1. Data (emails, texts)
2. Labels for each data point (spam/not spam, which author it was written by)
3. A way of “featurizing” the data that’s conducive to discriminating the classes
4. To know that it works.

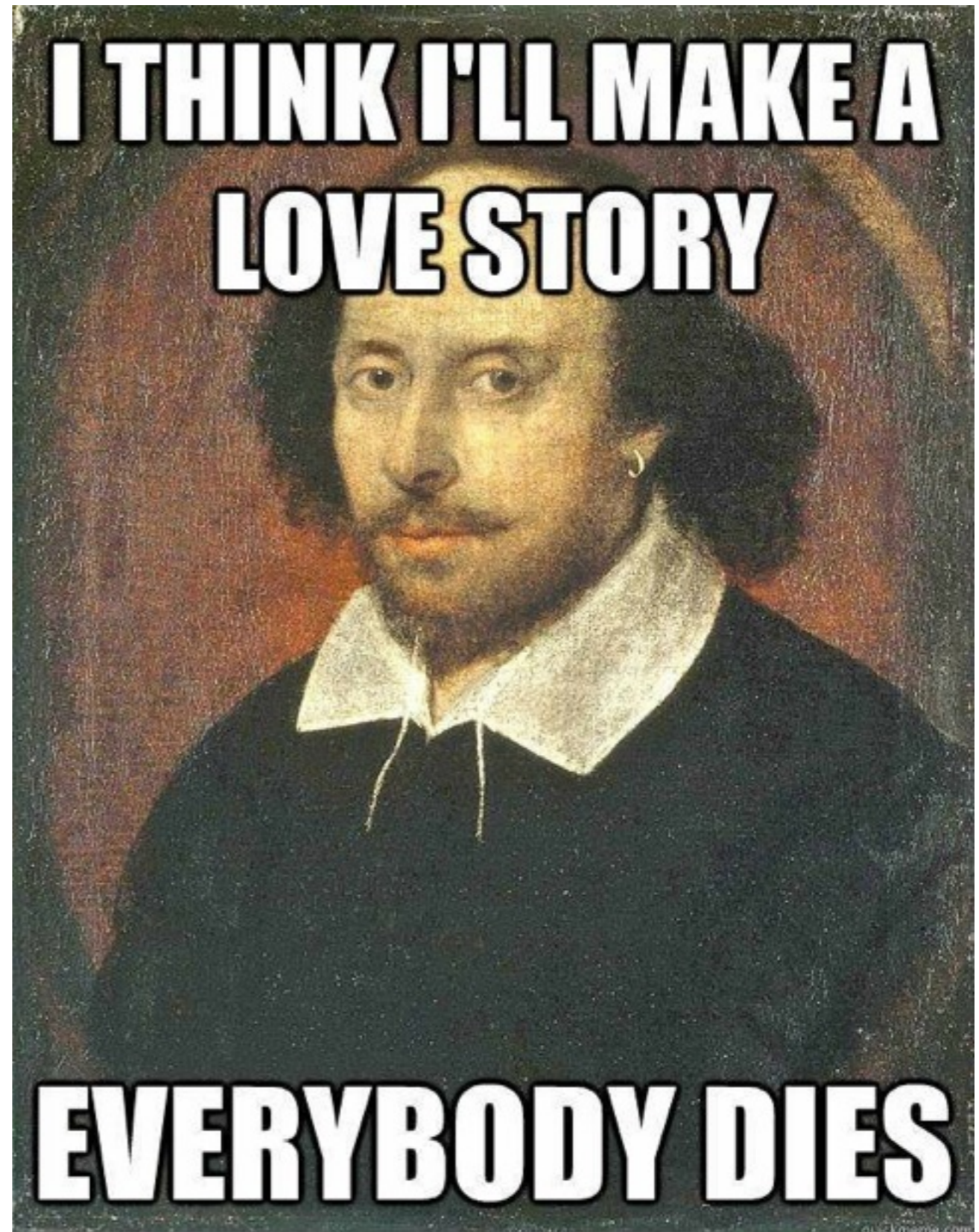
Recognizing a Classification Problem

- Can you formulate your question as a *choice* among some universe of possible classes?
- Can you create (or find) labeled data that marks that choice for a bunch of examples? Can ***you*** make that choice?
- Can you create features that might help in distinguishing those classes?

1. Those that belong to the emperor
2. Embalmed ones
3. Those that are trained
4. Suckling pigs
5. Mermaids (or Sirens)
6. Fabulous ones
7. Stray dogs
8. Those that are included in this classification
9. Those that tremble as if they were mad
10. Innumerable ones
11. Those drawn with a very fine camel hair brush
12. Et cetera
13. Those that have just broken the flower vase
14. Those that, at a distance, resemble flies

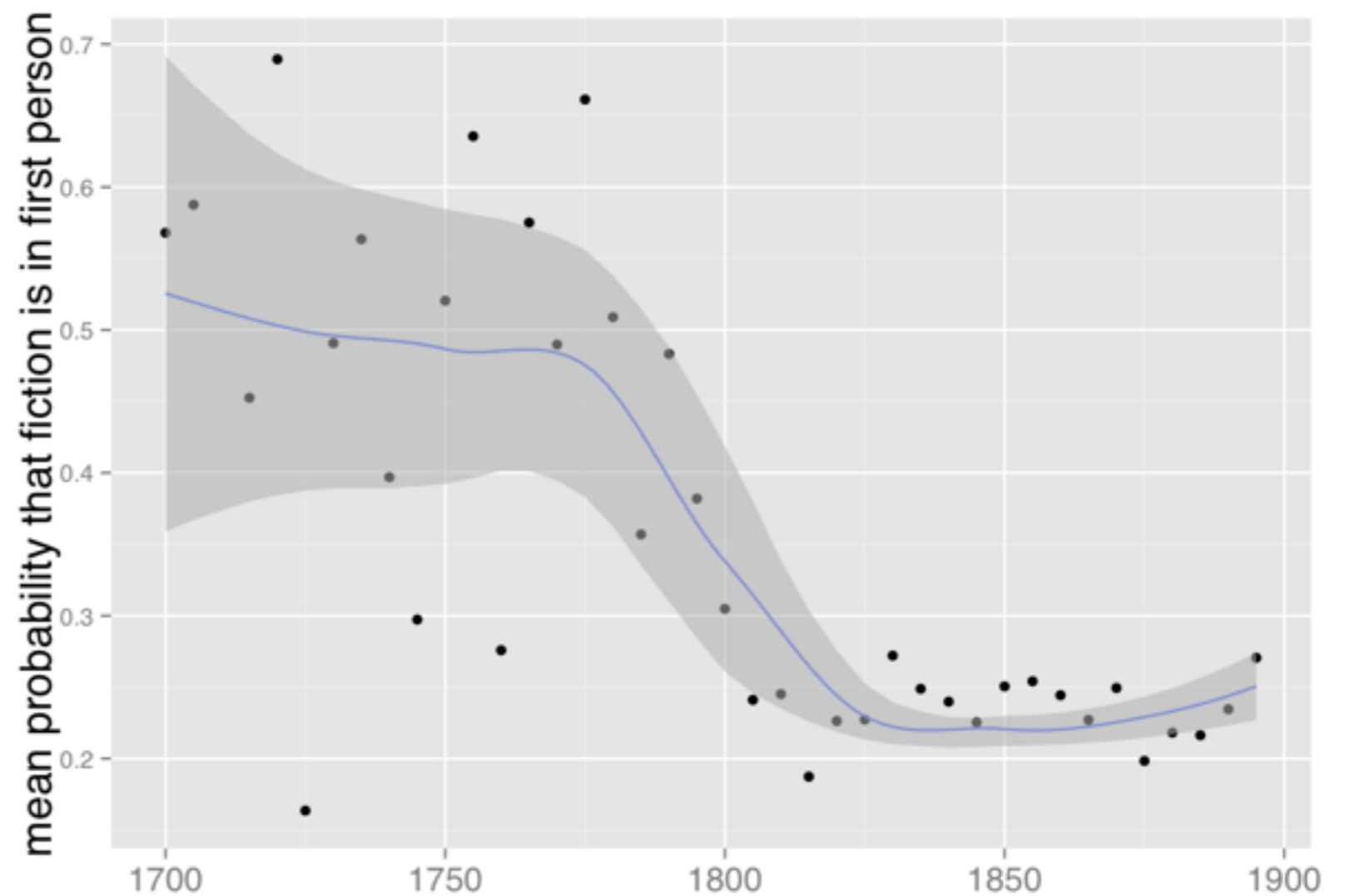


{Tragedy, Comedy}



Point of view

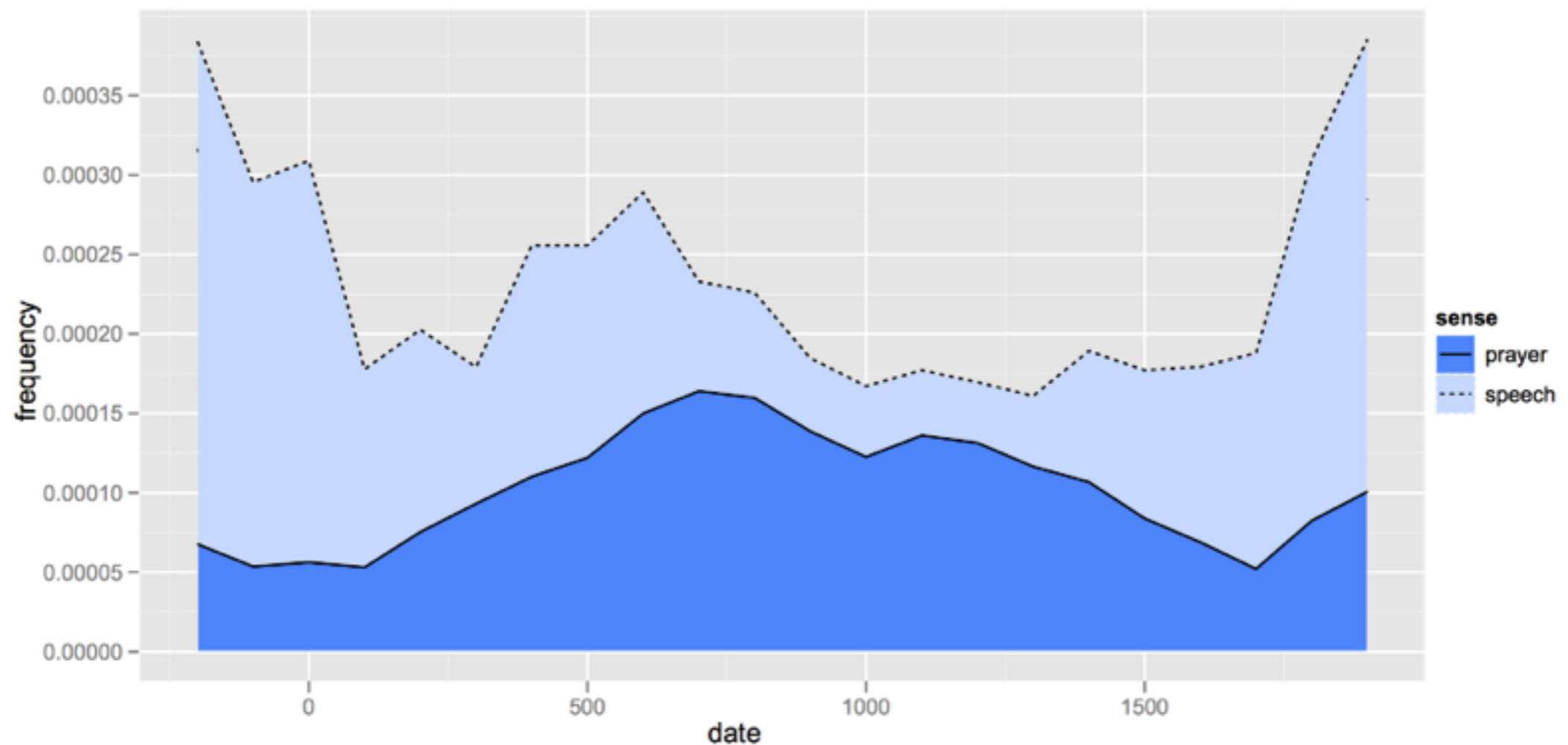
Classifying 1st-
vs. 3rd- person
narration in 32K
works of English-
language fiction



Ted Underwood, "Genre, gender and point of view" <http://tedunderwood.com/2013/09/22/genre-gender-and-point-of-view/>

Word sense

Classifying Latin “oratio” as *speech* vs. *prayer*.



Bamman and Crane, “Measuring Historical Word Sense Variation (JCDL 2011)”

Recognizing a Classification Problem

- I want to find all of the texts that have allusions to *Paradise Lost*.
- I want to know when discussions of “electricity” changed from magical to scientific.
- I want to find all of the “love oaths” in Shakespeare.

Classification Algorithms

- Naive Bayes
- Logistic Regression
- Support Vector Machines (SVM)
- Decision Trees/Random Forests
- K-nearest neighbors
- Hidden Markov Models (HMM)
- Conditional Random Fields (CRF)
- Structural SVM

Probability

- Lots of methods in the digital humanities/machine learning are *probabilistic*:
 - clustering, topic models
 - classification

Probability distributions

Normal

Gamma

Poisson

Geometric

Exponential

Multinomial

Bernoulli

Beta

Binomial

Uniform

Dirichlet

Probability distributions

Normal

Gamma

Poisson

Geometric

Exponential

Multinomial

Bernoulli

Beta

Binomial

Uniform

Dirichlet

Random variable

- A variable that can take values within a fixed set (discrete) or within some range (continuous).

$$X \in \{1, 2, 3, 4, 5, 6\}$$

$$X \in \{the, a, dog, cat, runs, to, store\}$$

$$P(X = x)$$

Probability that the random variable X takes the value x (e.g., 1)

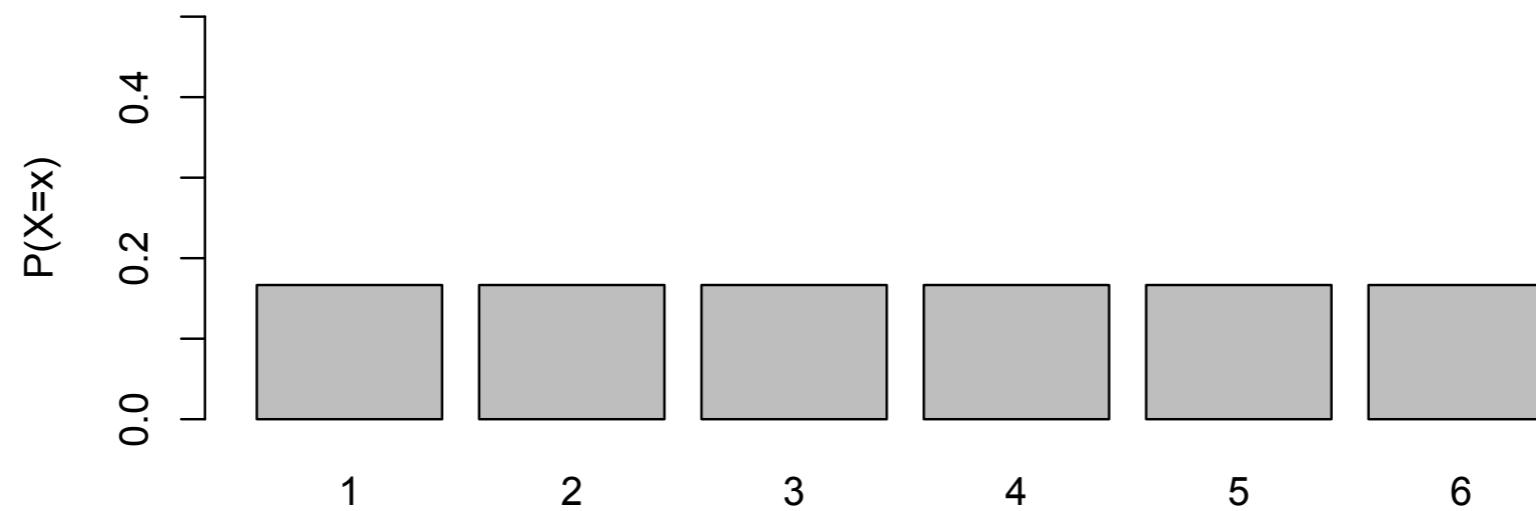
$$X \in \{1, 2, 3, 4, 5, 6\}$$

Two conditions:

1. Between 0 and 1: $0 \leq P(X = x) \leq 1$
2. Sum of all probabilities = 1 $\sum_x P(X = x) = 1$

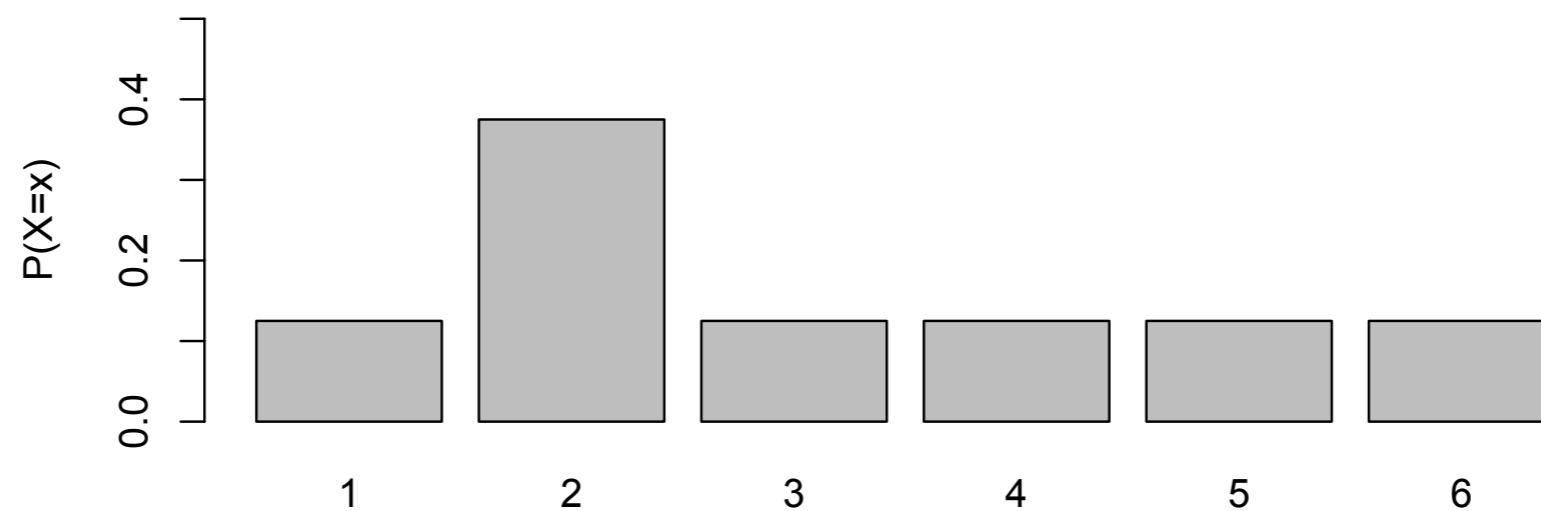
Fair dice

$$X \in \{1, 2, 3, 4, 5, 6\}$$



Weighted dice

$$X \in \{1, 2, 3, 4, 5, 6\}$$

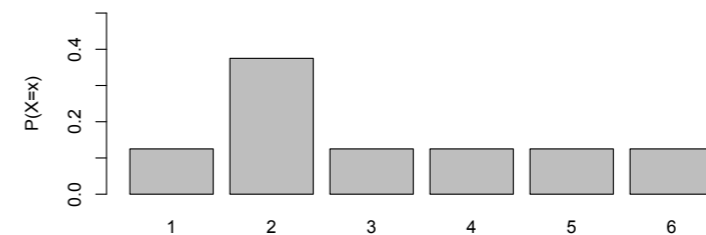
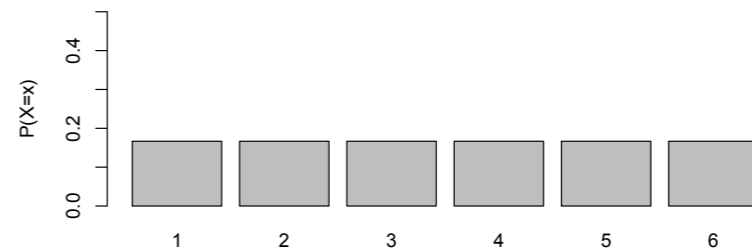
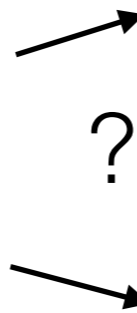


Parameter estimation

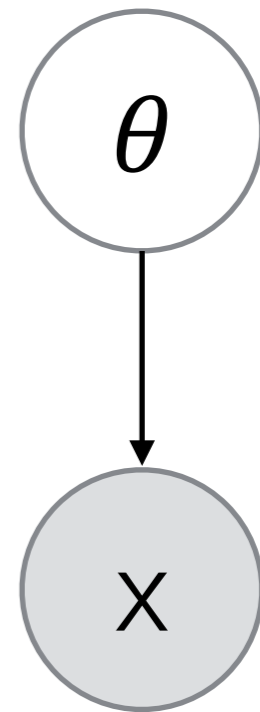
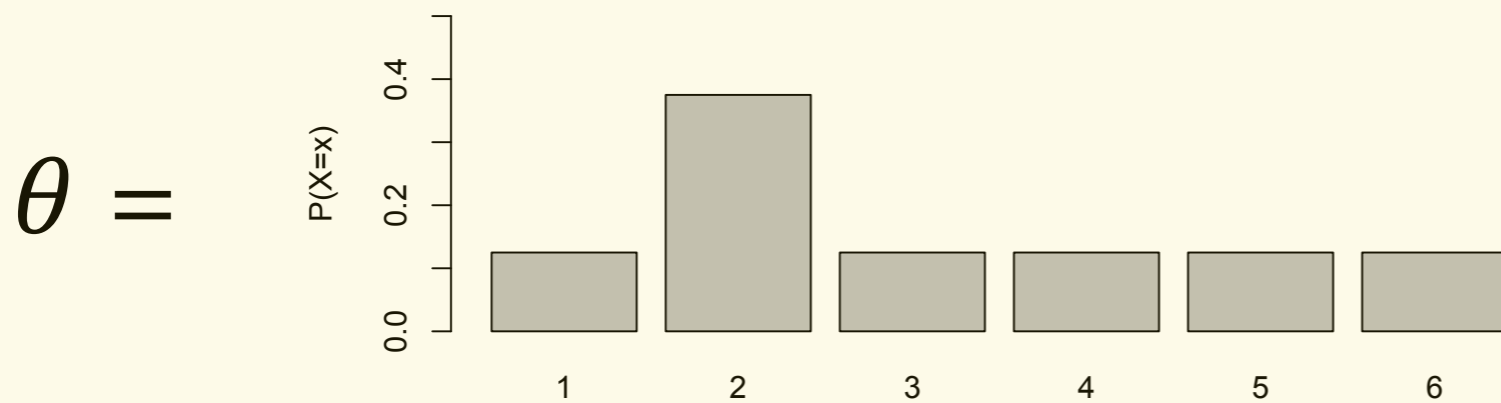
$$X \in \{1, 2, 3, 4, 5, 6\}$$

Data = 4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

We want to *estimate* the probability distribution that generated the data we see.



Generative story

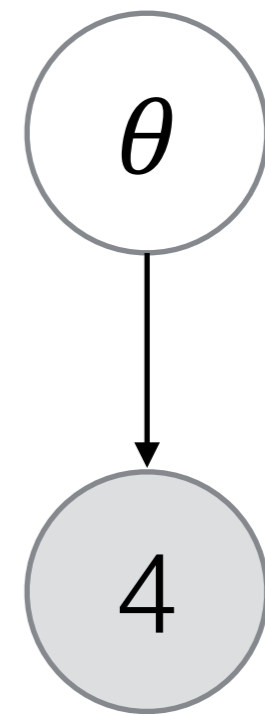
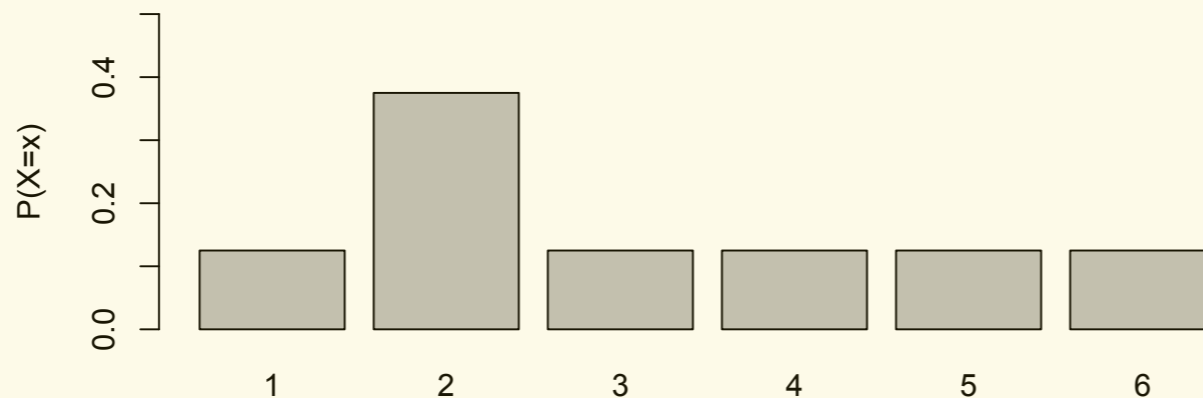


Data we see:

4,5,4,2,2,1,2,6,3,2,2,2,1,4,2

Generative story

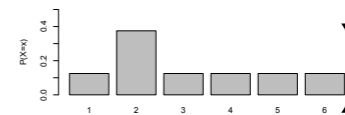
$\theta =$



Data we see:

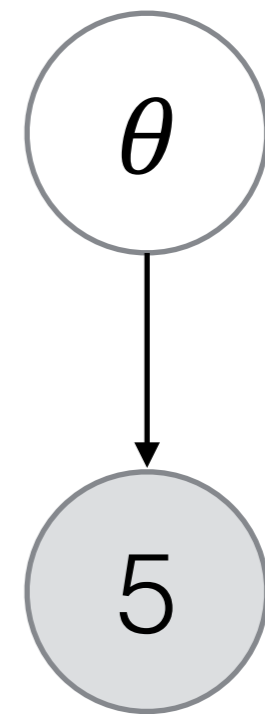
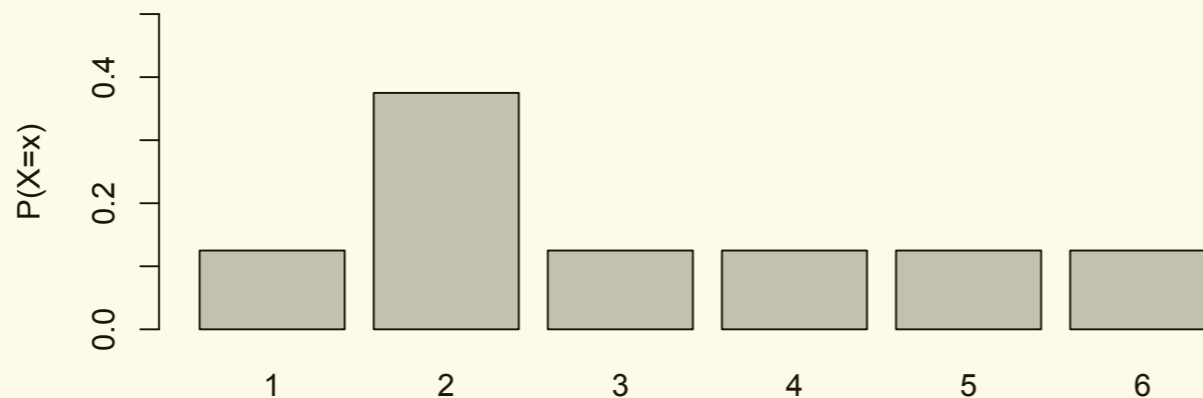
4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

$$P(X=4|\theta = \text{distribution}) = .125$$



Generative story

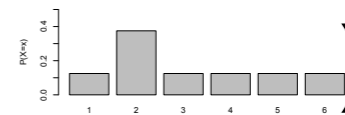
$\theta =$



Data we see:

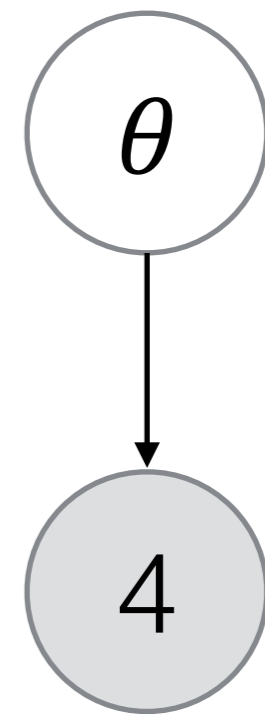
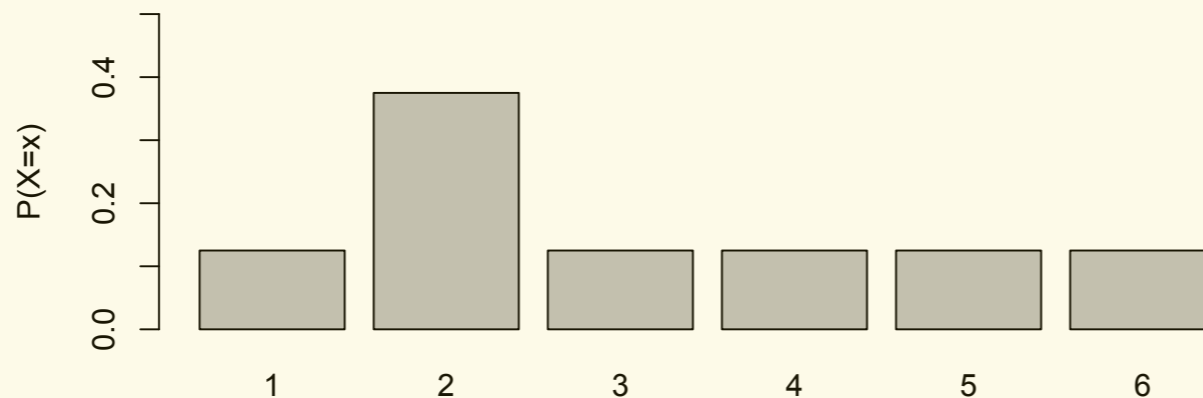
4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

$$P(X=5|\theta = \text{distribution}) = .125$$



Generative story

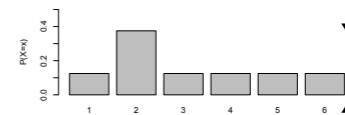
$\theta =$



Data we see:

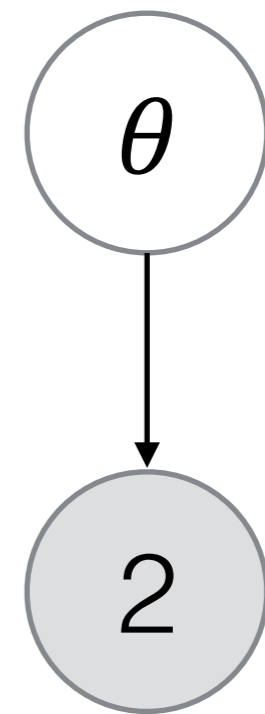
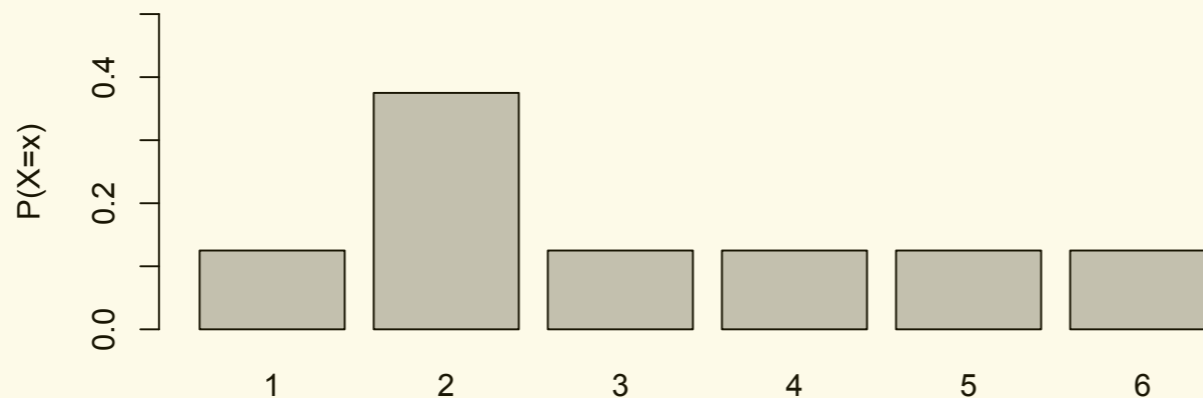
4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

$$P(X=4|\theta = \text{distribution}) = .125$$



Generative story

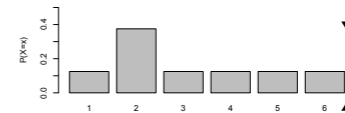
$\theta =$



Data we see:

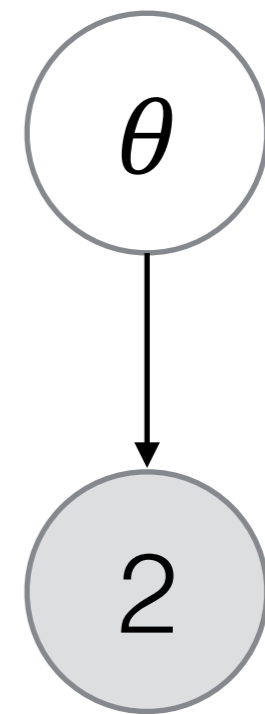
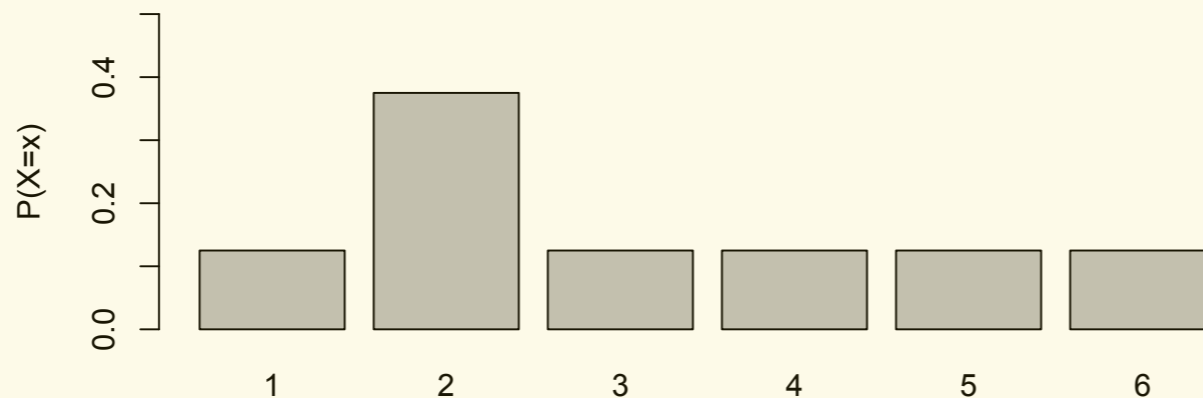
4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

$$P(X=2|\theta=) = .375$$



Generative story

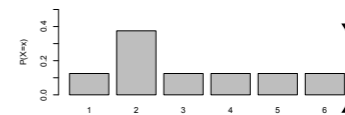
$\theta =$



Data we see:

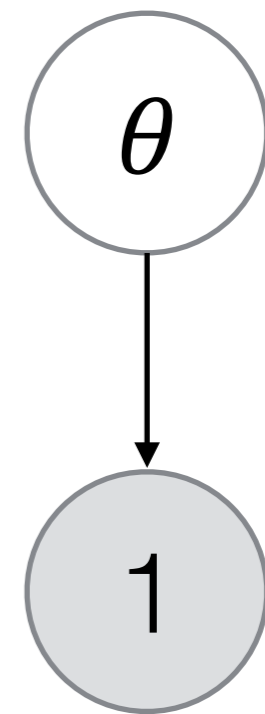
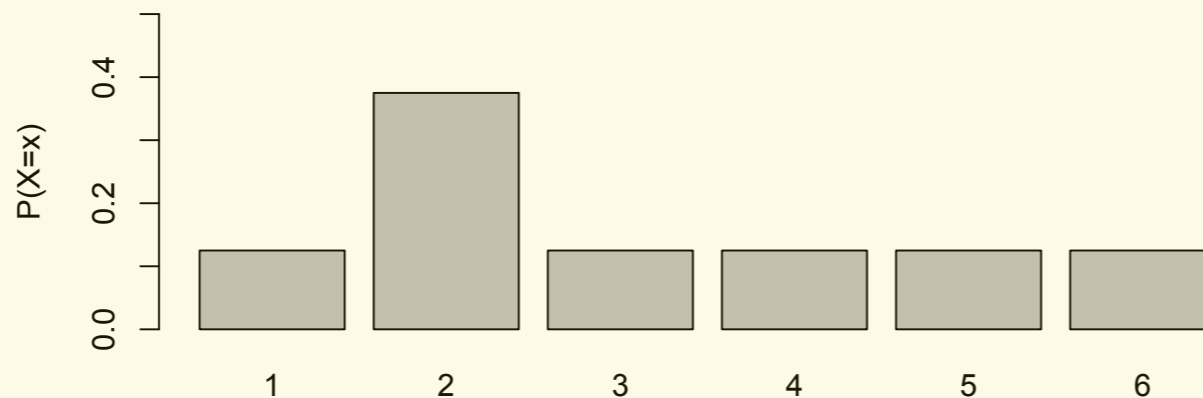
4, 5, 4, 2, 2, 1, 2, 6, 3, 2, 2, 2, 1, 4, 2

$$P(X=2|\theta = \text{distribution}) = .375$$



Generative story

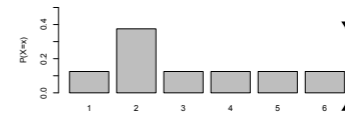
$\theta =$



Data we see:

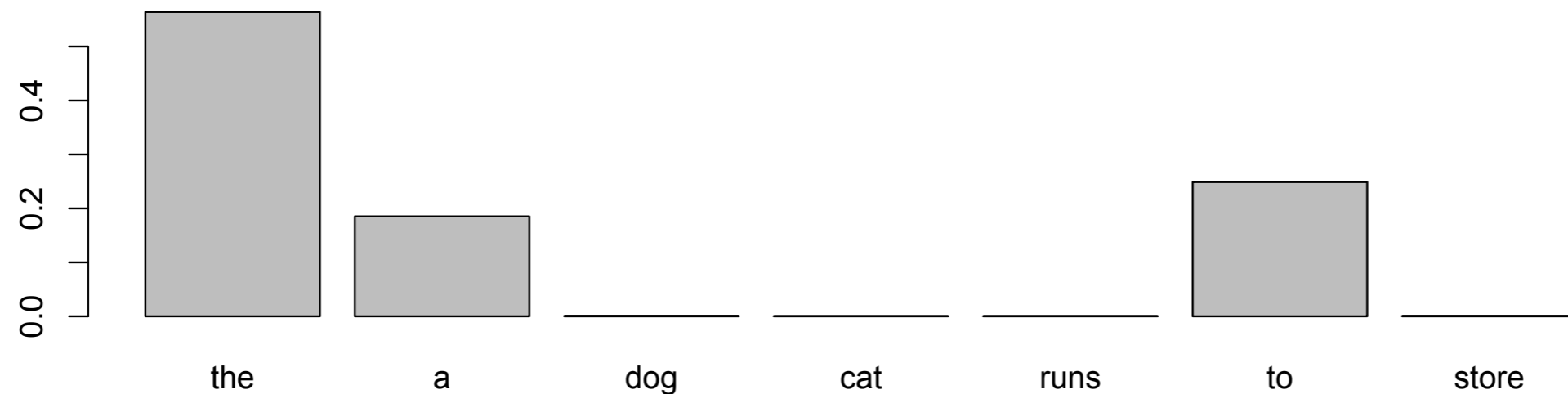
4,5,4,2,2,1,2,6,3,2,2,2,1,4,2

$$P(X=1|\theta = \text{distribution}) = .125$$



Unigram probability

$$X \in \{the, a, dog, cat, runs, to, store\}$$



How do we calculate this?

$$P(X=\text{"the"}) = 28/536 = .052$$

Conditional Probability

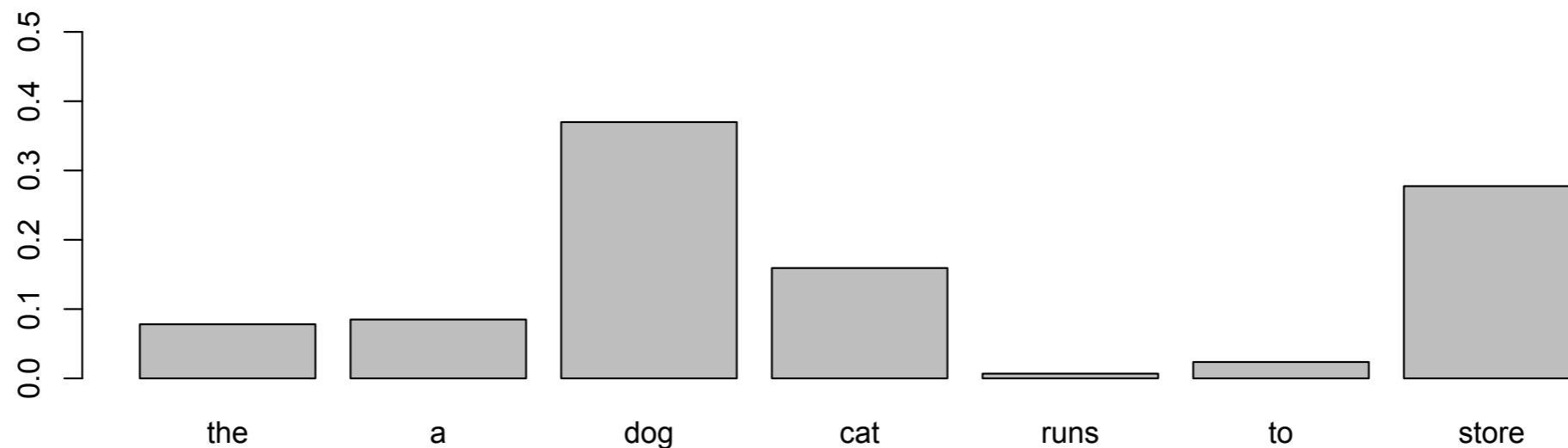
$$P(X = x|Y = y)$$

- Probability that one random variable takes a particular value *given* the fact that a different variable takes another

$$P(X_i = \text{dog}|X_{i-1} = \text{the})$$

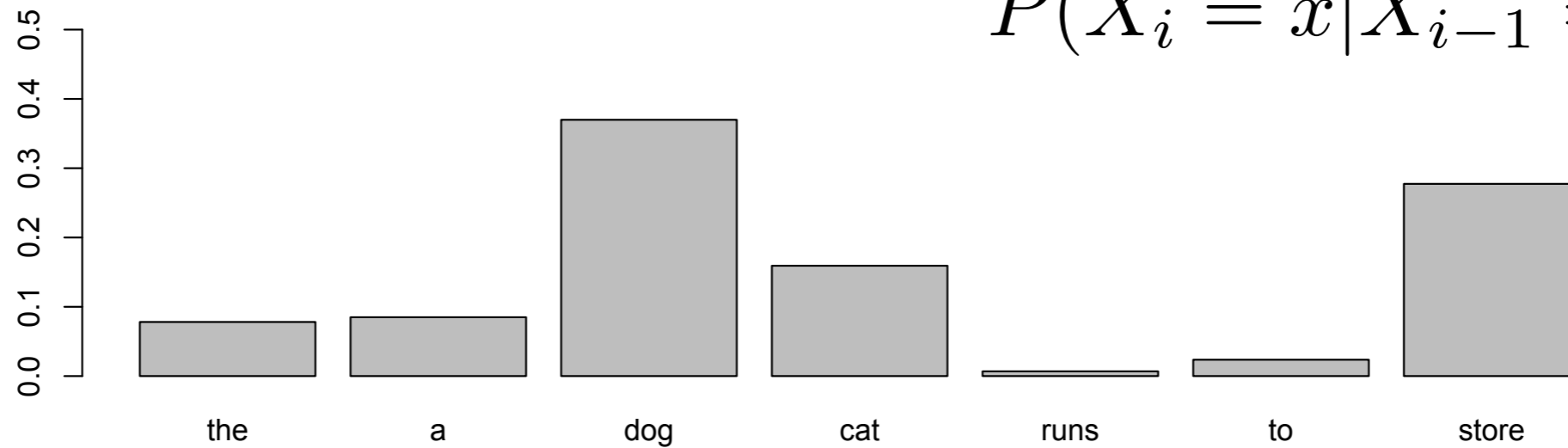
Conditional Probability

$$P(X_i = \text{dog} | X_{i-1} = \text{the})$$

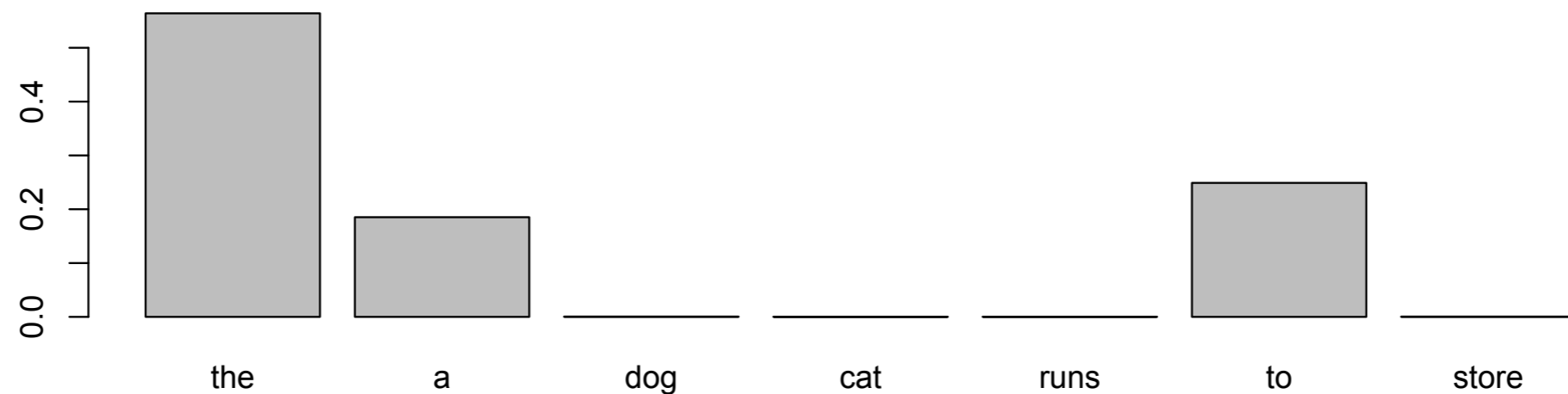


Conditional Probability

$$P(X_i = x | X_{i-1} = the)$$



$$P(X_i = x)$$



$$P(X_i = \text{"room"} | X_{i-1} = \text{"the"}) = 2/28 = .071$$

Conditional Probability

$P(X = \text{vampire})$ vs. $P(X = \text{vampire} | Y = \text{horror})$

$P(X = \text{manners} | Y = \text{austen})$ vs. $P(X = \text{whale} | Y = \text{austen})$

$P(X = \text{manners} | Y = \text{austen})$ vs. $P(X = \text{manners} | Y = \text{dickens})$

Our first classifier

“Mr. Collins was not a sensible man”

Austen		Dickens	
$P(X=\text{Mr.} \mid Y=\text{Austen})$	0.0084	$P(X=\text{Mr.} \mid Y=\text{Dickens})$	0.00421
$P(X=\text{Collins} \mid Y=\text{Austen})$	0.00036	$P(X=\text{Collins} \mid Y=\text{Dickens})$	0.000016
$P(X=\text{was} \mid Y=\text{Austen})$	0.01475	$P(X=\text{was} \mid Y=\text{Dickens})$	0.015043
$P(X=\text{not} \mid Y=\text{Austen})$	0.01145	$P(X=\text{not} \mid Y=\text{Dickens})$	0.00547
$P(X=\text{a} \mid Y=\text{Austen})$	0.01591	$P(X=\text{a} \mid Y=\text{Dickens})$	0.02156
$P(X=\text{sensible} \mid Y=\text{Austen})$	0.00025	$P(X=\text{sensible} \mid Y=\text{Dickens})$	0.00005
$P(X=\text{man} \mid Y=\text{Austen})$	0.00121	$P(X=\text{man} \mid Y=\text{Dickens})$	0.001707

Our first classifier

“Mr. Collins was not a sensible man”

$P(X = \text{“Mr. Collins was not a sensible man”} \mid Y = \text{Austen})$

$$\begin{aligned} &= P(\text{“Mr”} \mid \text{Austen}) \times P(\text{“Collins”} \mid \text{Austen}) \times \\ &P(\text{“was”} \mid \text{Austen}) \times P(\text{“not”} \mid \text{Austen}) \dots \\ &= 0.000000022507322 \ (\approx \mathbf{2.3 \times 10^{-8}}) \end{aligned}$$

$P(X = \text{“Mr. Collins was not a sensible man”} \mid Y = \text{Dickens})$

$$\begin{aligned} &P(\text{“Mr”} \mid \text{Dickens}) \times P(\text{“Collins”} \mid \text{Dickens}) \times \\ &P(\text{“was”} \mid \text{Dickens}) \times P(\text{“not”} \mid \text{Dickens}) \dots \\ &= 0.000000002078906 \ (\approx \mathbf{2.1 \times 10^{-9}}) \end{aligned}$$

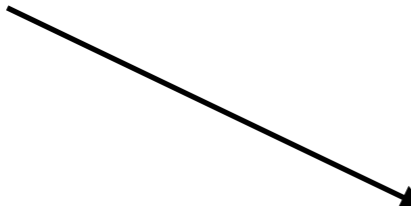
Bayes' Rule

$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

Bayes' Rule

Prior belief that $Y = y$
(before you see any data)

Likelihood of the data
given that $Y=y$



$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

Posterior belief that $Y=y$ given that $X=x$

Bayes' Rule

Prior belief that $Y = y$
(before you see any data)

Likelihood of the data
given that $Y=y$

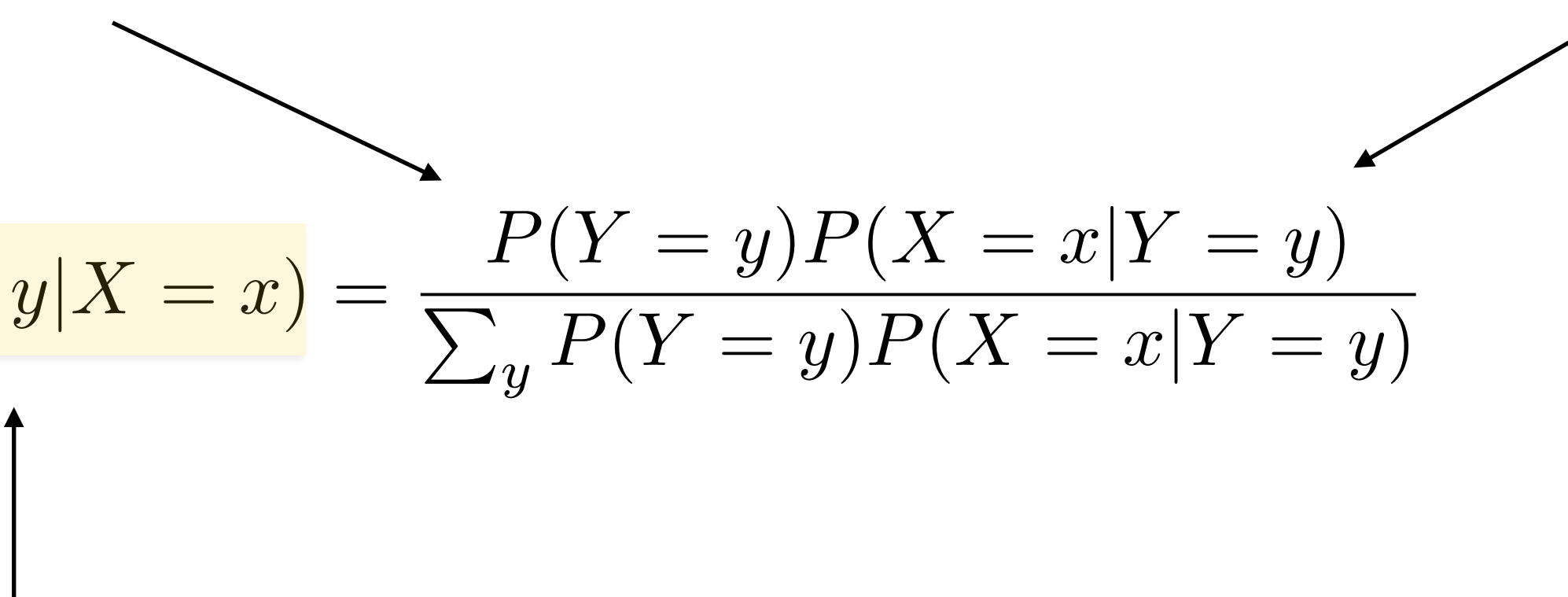

$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

Posterior belief that $Y=y$ given that $X=x$

Bayes' Rule

Prior belief that $Y = y$
(before you see any data)

Likelihood of the data
given that $Y=y$


$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

Posterior belief that $Y=y$ given that $X=x$

Bayes' Rule

Prior belief that $Y = \text{Austen}$
(before you see any data)

Likelihood of “Mr. Collins
was not a sensible man”
given that $Y = \text{Austen}$

$$P(Y = y | X = x) = \frac{P(Y = y)P(X = x | Y = y)}{\sum_y P(Y = y)P(X = x | Y = y)}$$

Posterior belief that $Y = \text{Austen}$ given that
 $X = \text{“Mr. Collins was not a sensible man”}$

This sum ranges over
 $y = \text{Austen} + y = \text{Dickens}$
(so that it sums to 1)

Naive Bayes Classifier

$$\frac{P(Y = \textit{Austen})P(X = \textit{“Mr...”} | Y = \textit{Austen})}{P(Y = \textit{Austen})P(X = \textit{“Mr...”} | Y = \textit{Austen}) + P(Y = \textit{Dickens})P(X = \textit{“Mr...”} | Y = \textit{Dickens})}$$

Let's say $P(Y=\textit{Austen}) = P(Y=\textit{Dickens}) = 0.5$
(i.e., both are equally likely a priori)

$$= \frac{0.5 \times (2.3 \times 10^{-8})}{0.5 \times (2.3 \times 10^{-8}) + 0.5 \times (2.1 \times 10^{-9})}$$

$$P(Y = \textit{Austen} | X = \textit{“Mr...”}) = 91.5\%$$

$$P(Y = \textit{Dickens} | X = \textit{“Mr...”}) = 8.5\%$$

Taxicab Problem

“A cab was involved in a hit and run accident at night. Two cab companies, the Green and the Blue, operate in the city. You are given the following data:

- 85% of the cabs in the city are Blue.
- A witness identified the cab as Blue. The witness was interviewed on the night of the accident and concluded that the witness correctly identified each one of the two colors 80% of the time and failed 20% of the time.

“Base rate fallacy”
Don't ignore prior information!

What is the probability that the cab involved in the accident was Blue rather than Green knowing that this witness identified it as Blue?”

(Tversky & Kahneman 1981)

Prior Belief

- Now let's assume that Dickens published 1000 times more books than Austen.
- $P(Y = \text{Austen}) = 0.000999$
- $P(Y = \text{Dickens}) = 0.999001$

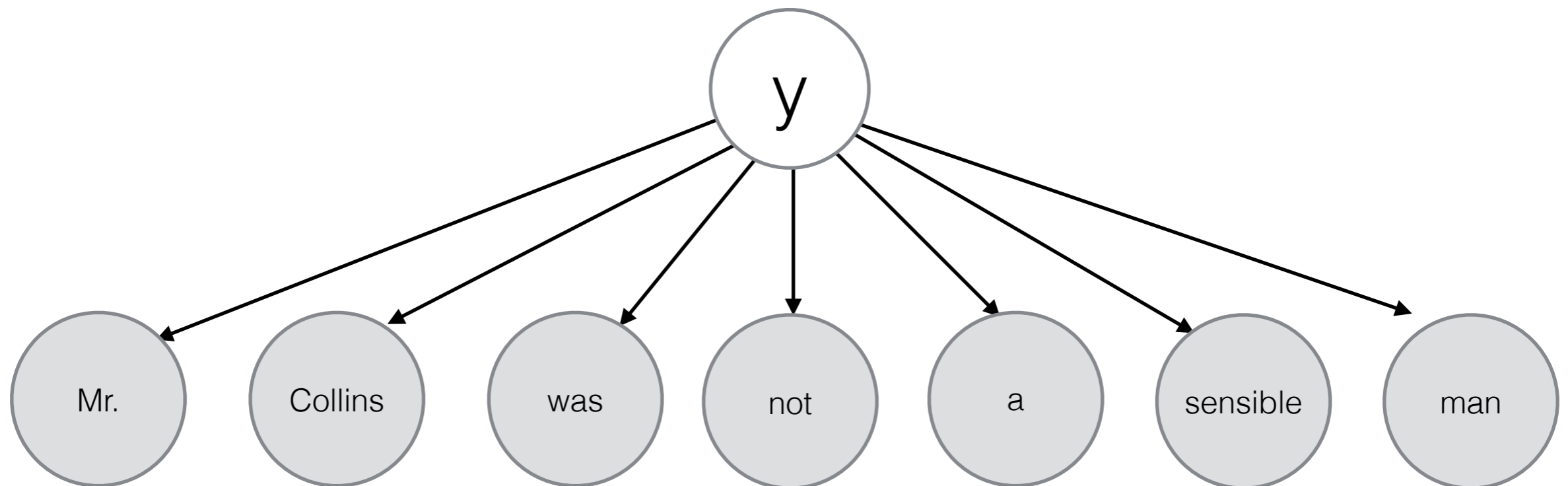
$$\frac{0.000999 \times (2.3 \times 10^{-8})}{0.000999 \times (2.3 \times 10^{-8}) + 0.999001 \times (2.1 \times 10^{-9})}$$

$$P(Y = \text{Austen} | X) = 0.011$$

$$P(Y = \text{Dickens} | X) = 0.989$$

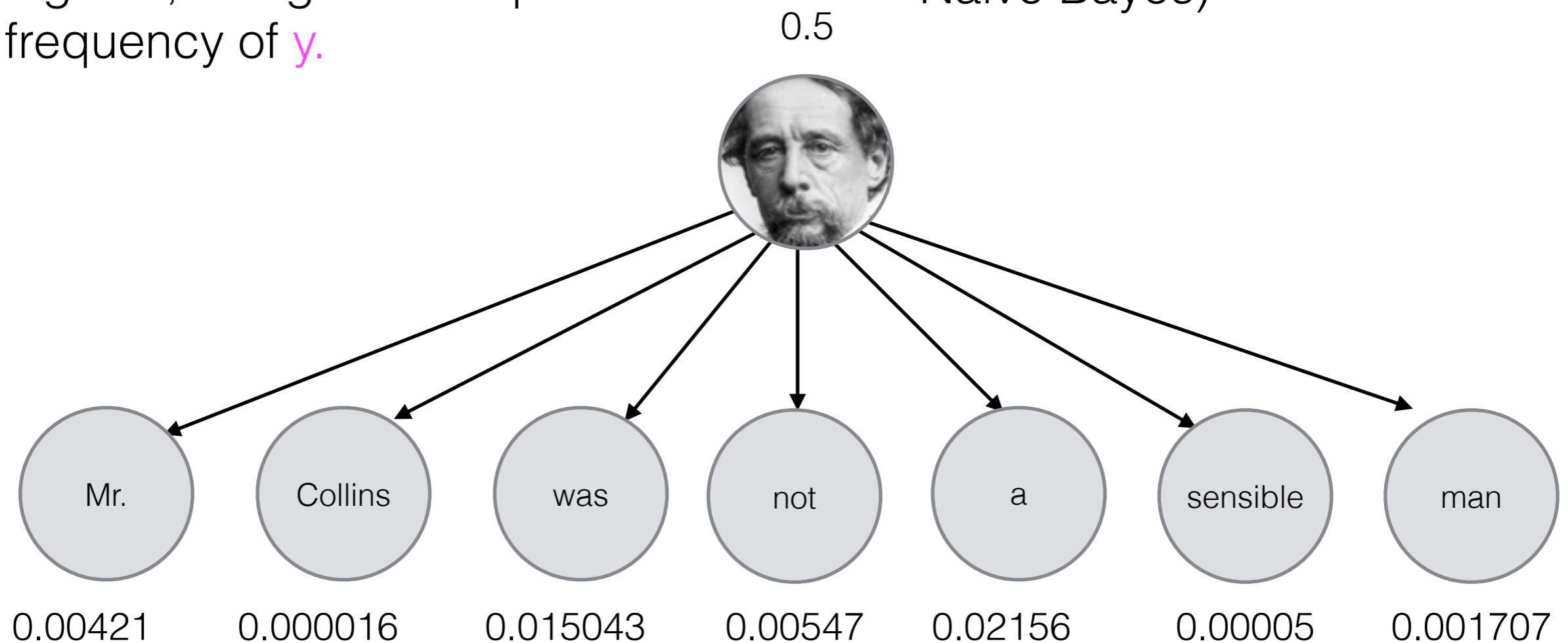
Naive Bayes

- Find the value of y (e.g., author) for which the probability of the x (e.g., the words) that we see is highest, along with the prior frequency of y .
- All x 's are independent and contribute equally to finding the best y (the “naive” in Naive Bayes)



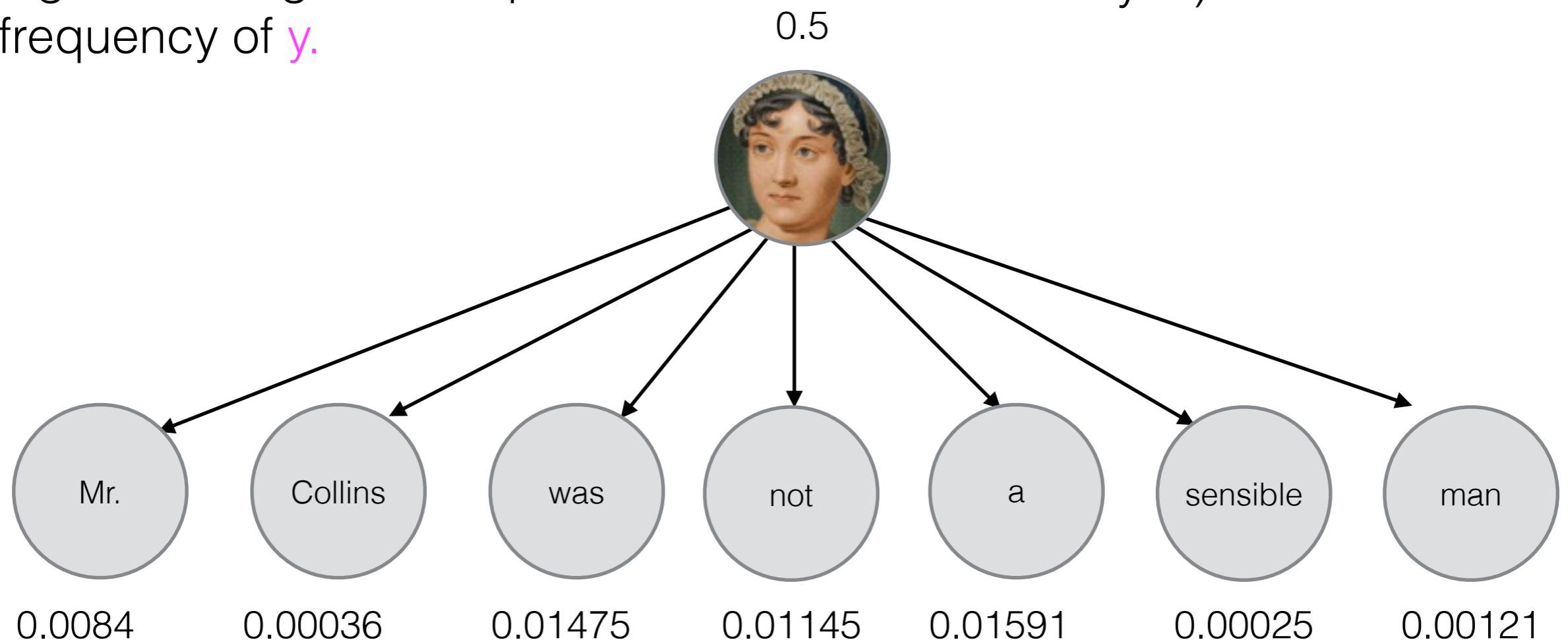
Naive Bayes

- Find the value of y (e.g., author) for which the probability of the x (e.g., the words) that we see is highest, along with the prior frequency of y .
- All x 's are independent and contribute equally to finding the best y (the “naive” in Naive Bayes)



Naive Bayes

- Find the value of y (e.g., author) for which the probability of the x (e.g., the words) that we see is highest, along with the prior frequency of y .
- All x 's are independent and contribute equally to finding the best y (the “naive” in Naive Bayes)



Parameters

$P(X = x \mid Y = \text{Austen})$

value	prob
Mr.	0.0084
Collins	0.00036
was	0.01475
not	0.01145
a	0.01591
sensible	0.00025
man	0.00121
dog	0.003
chimney	0.004

...

$P(X = x \mid Y = \text{Dickens})$

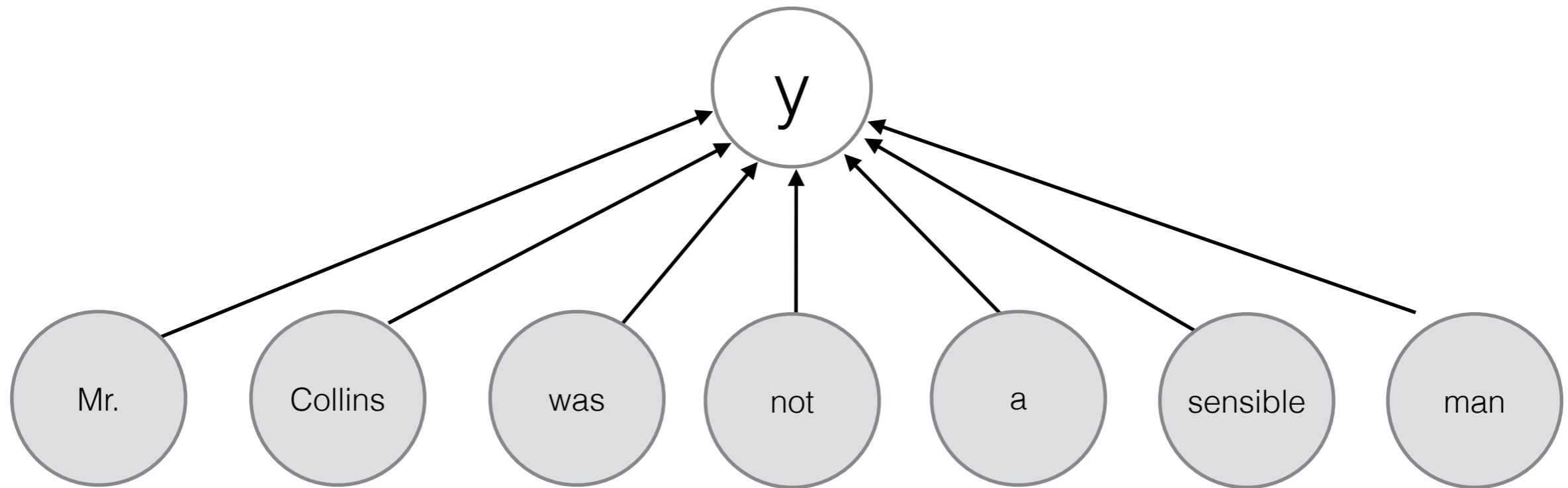
value	prob
Mr.	0.00421
Collins	0.000016
was	0.015043
not	0.00547
a	0.02156
sensible	0.00005
man	0.001707
dog	0.002
chimney	0.008

...

$P(Y = y)$

value	prob
Dickens	0.50
Austen	0.50

Logistic Regression



$$P(Y = y|X, \beta) = \frac{\exp \left(\sum_i^F \beta_{y,i} x_i \right)}{\sum_{y'} \exp \left(\sum_i^F \beta_{y',i} x_i \right)}$$

Logistic Regression

$\beta_{\text{Austen}} =$

i	feat	value
1	Mr.	1.4
2	Collins	15.7
3	was	0.01
4	a	-0.003
5	sensible	7.8
6	man	1.3
7	dog	-1.3
8	chimney	-10.3

$X =$

i	feat	value
1	Mr.	1
2	Collins	1
3	was	1
4	a	1
5	sensible	1
6	man	1
7	dog	0
8	chimney	0

$$P(Y = y|X, \beta) = \frac{\exp \left(\sum_i^F \beta_{y,i} x_i \right)}{\sum_{y'} \exp \left(\sum_i^F \beta_{y',i} x_i \right)}$$

Logistic Regression

- Find the value of β that maximizes $P(Y=y \mid X=x, \beta)$ where we know the value of y given a particular x (i.e., in training data).

- Likelihood:
$$L(\beta) = \prod_{\{x,y\}} P(Y = y \mid X = x, \beta)$$

- Log Likelihood:
$$\ell(\beta) = \sum_{\{x,y\}} \log P(Y = y \mid X = x, \beta)$$

Overfitting

- Memorizing patterns in the training data **too well** → perform worse on data you don't train on.
- e.g., if we see **Collins** only in Austen books in the training data, what happens if we see **Collins** in a new book we're predicting?

Regularization

- Penalize parameters that are very big (i.e., that are far away from 0).

i	feat	β
1	Mr.	1.4
2	Collins	18403.0
3	was	0.01
4	a	-0.003
5	sensible	7.8
6	man	1.3
7	dog	-1.3
8	chimney	-10.3

$$\arg \max_{\beta} \sum_{\{x,y\}} \log P(Y = y | X = x, \beta) - \lambda \sum_j \beta_j^2$$

Regularization

i	feat	β
1	Mr.	1.4
2	Collins	18403
3	was	0.01
4	a	-0.003
5	sensible	7.8
6	man	1.3
7	dog	-1.3
8	chimney	-10.3



i	feat	β
1	Mr.	1.1
2	Collins	13.8
3	was	0.005
4	a	-0.0007
5	sensible	6.9
6	man	0.9
7	dog	-0.7
8	chimney	-8.3

$$\arg \max_{\beta} \sum_{\{x,y\}} \log P(Y = y | X = x, \beta) - \lambda \sum_j \beta_j^2$$

Regularization

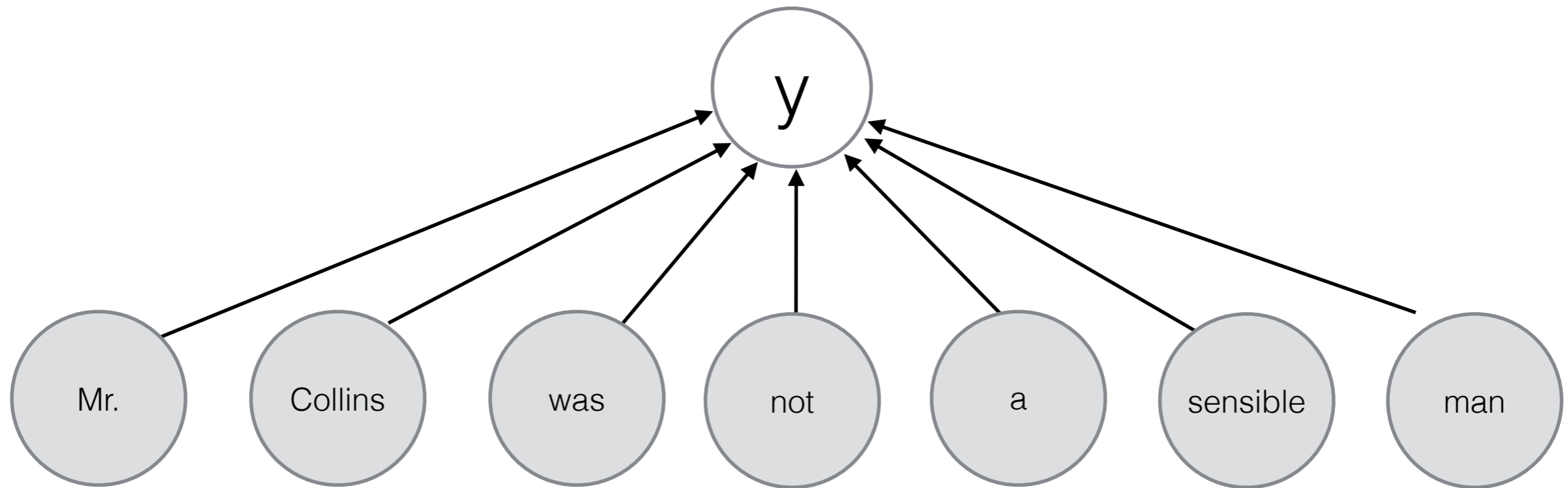
- L2 regularization encourages parameters to be close to 0.

$$\arg \max_{\beta} \sum_{\{x,y\}} \log P(Y = y|X = x, \beta) - \lambda \sum_j \beta_j^2$$

- L1 regularization also encourages them to be 0. (Sparsity)

$$\arg \max_{\beta} \sum_{\{x,y\}} \log P(Y = y|X = x, \beta) - \lambda \sum_j |\beta_j|$$

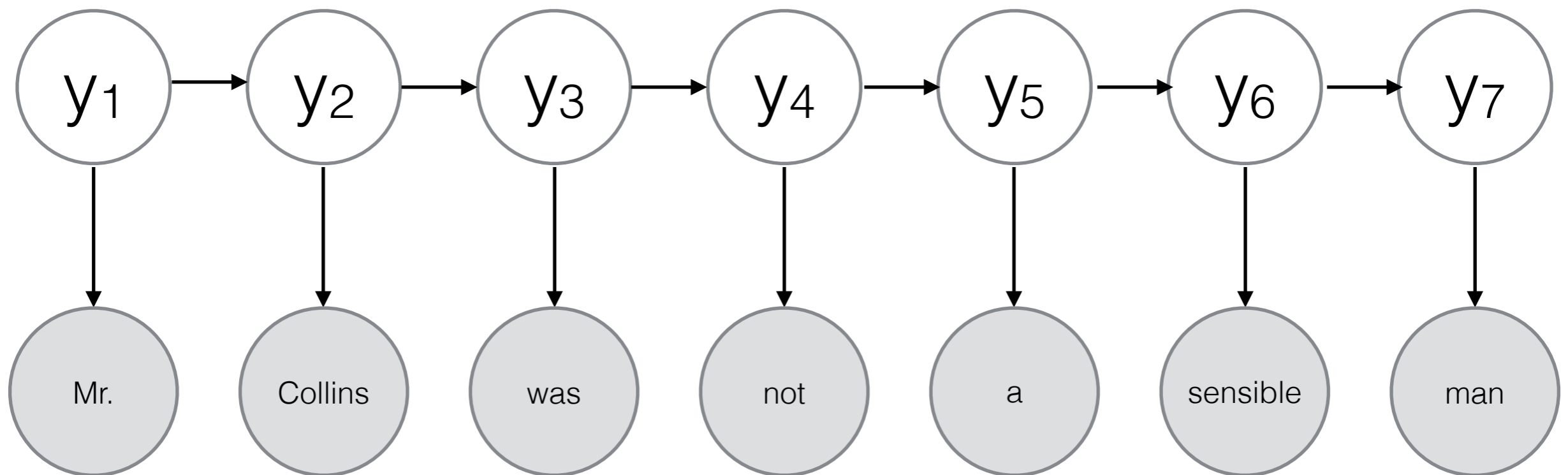
Logistic Regression



$$P(Y = y|X, \beta) = \frac{\exp \left(\sum_i^F \beta_{y,i} x_i \right)}{\sum_{y'} \exp \left(\sum_i^F \beta_{y',i} x_i \right)}$$

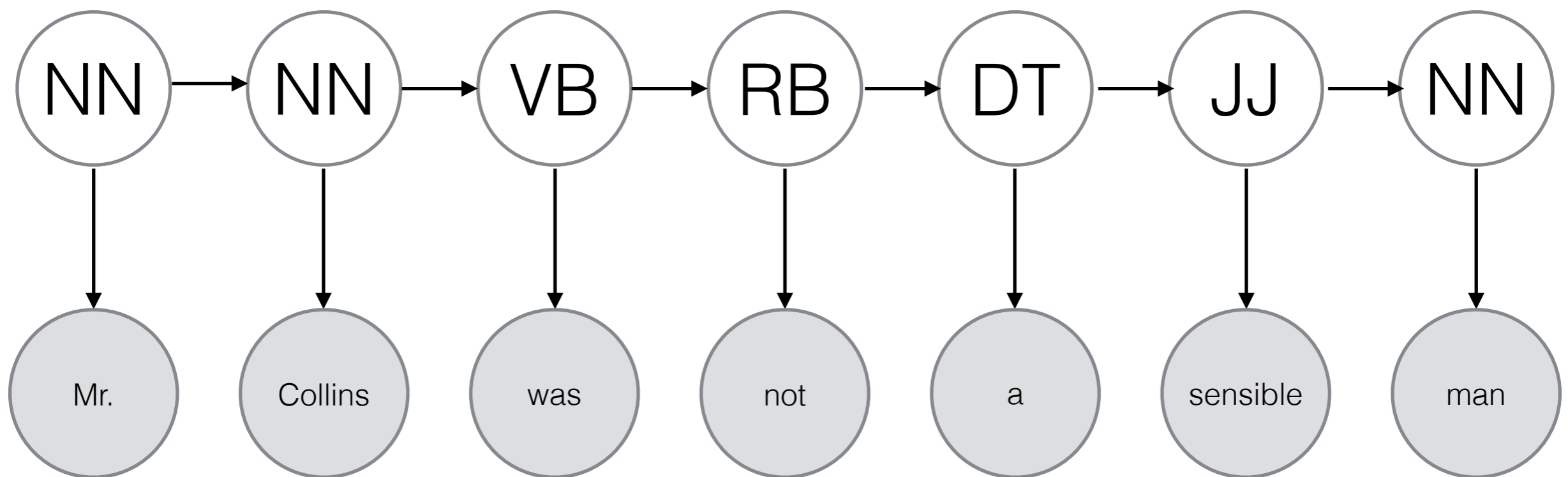
Hidden Markov Model

Generative model for predicting a sequence of variables.



Hidden Markov Model

Example: part of speech tagging



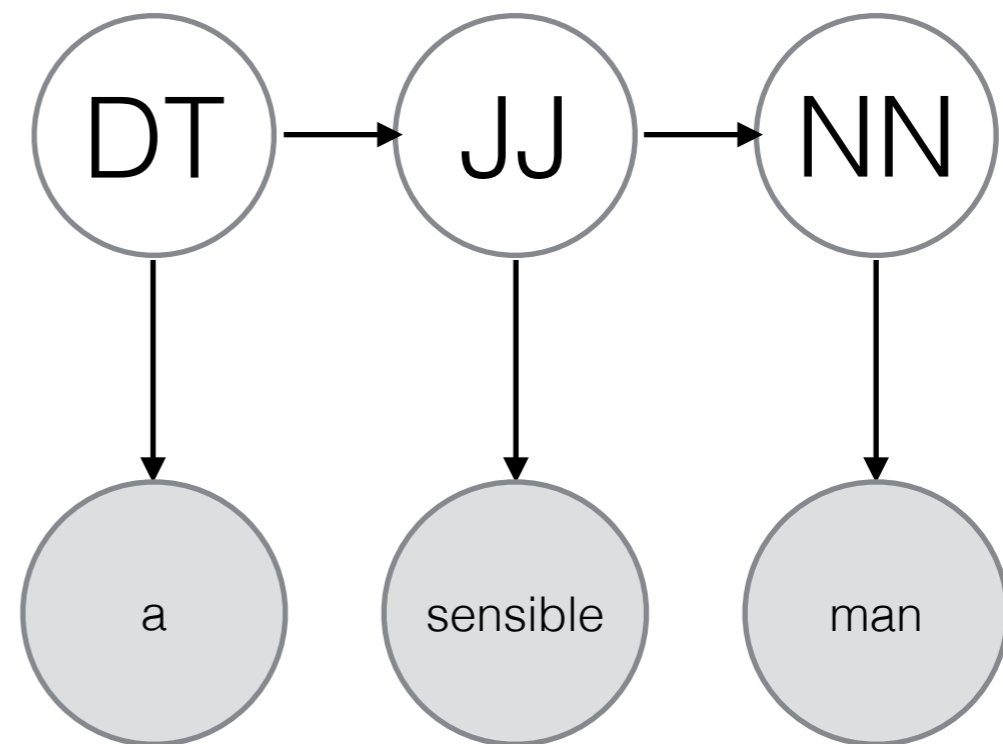
Hidden Markov Model

$$P(X=x \mid y = \text{DT})$$

value	prob
a	0.37
the	0.33
an	0.17
sensible	0
dog	0

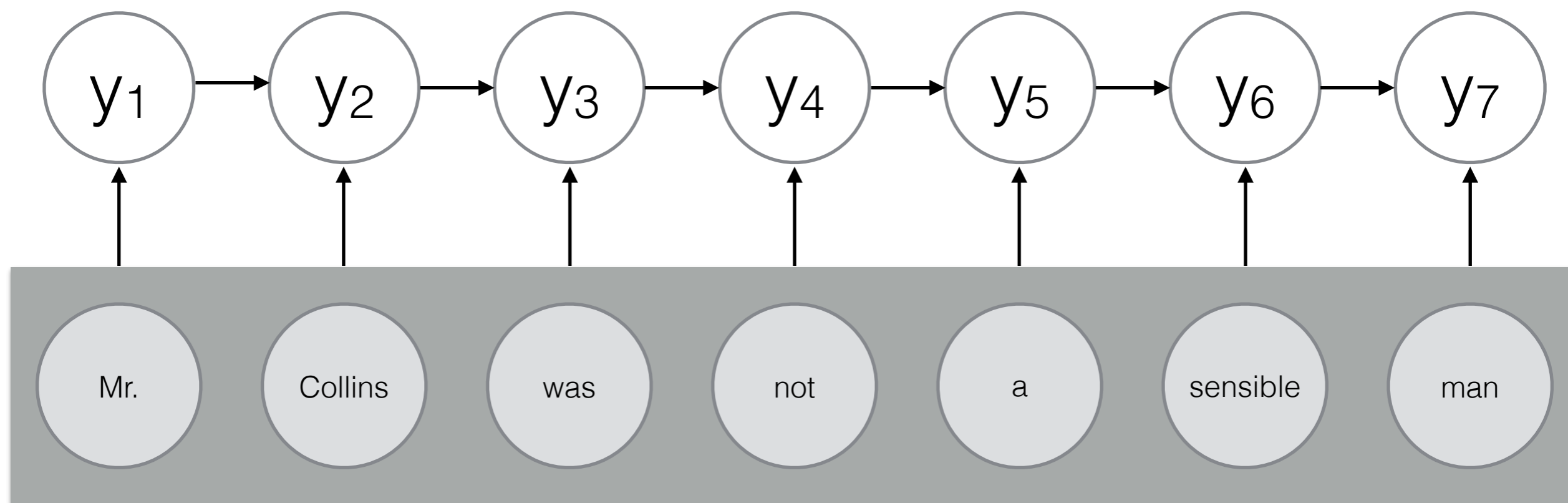
$$P(Y_i = y \mid Y_{i-1} = \text{DT})$$

value	prob
NN	0.38
JJ	0.17
RB	0.15
DT	0



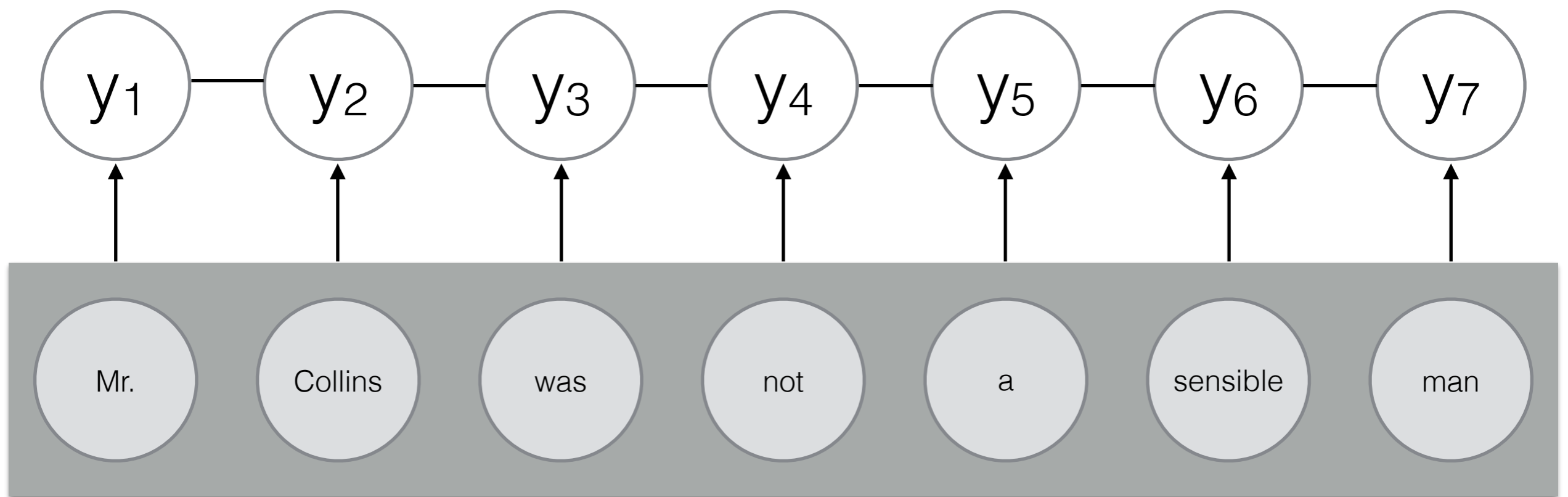
Maximum Entropy Markov Model

Discriminative model for predicting a sequence of variables.



Conditional Random Field

Discriminative model for predicting a sequence of variables.



Rich features

Mr. **Collins** was not a sensible man

HMM

feature	val
word=Collins	1
word=the	0
word=a	0
word=not	0
word=sensible	0

MEMM/CRF

feature	val
word=Collins	1
word starts with capital letter	1
word is in list of known names	1
word ends in -ly	0

Try it yourself

- LightSide
<http://bit.ly/1hdKX0R>
- (Google “LightSide Academic”)

Break!

Unsupervised Learning

- Unsupervised learning finds *interesting structure* in data.
- clustering data into groups
- discovering “factors”
- discovering graph structure (6DFB)



Unsupervised Learning

- Matrix completion (e.g., user recommendations on Netflix, Amazon)

	Ann	Bob	Chris	David	Erik
Star Wars	5	5	4	5	3
Bridget Jones		4		4	1
Rocky	3		5		
Rambo		?		2	5

Unsupervised Learning

- Hierarchical clustering
- Flat clustering (K-means)
- Topic models

Hierarchical Clustering

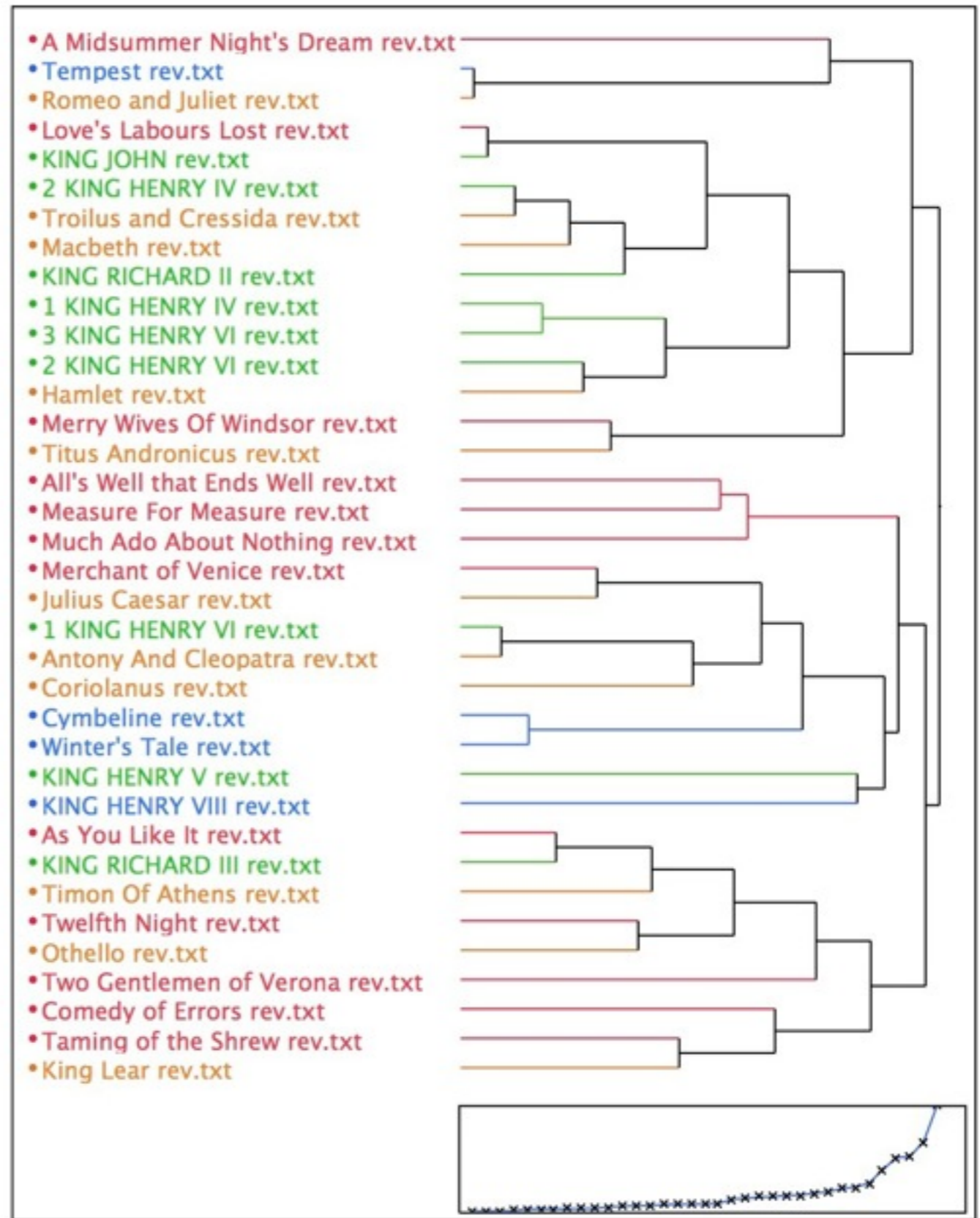
- *Hierarchical* order among the elements being clustered
- Bottom-up = agglomerative clustering
- Top-down = divisive clustering

Dendrogram

Shakespeare's plays

Witmore (2009)

<http://winedarksea.org/?p=519>



Bottom-up clustering

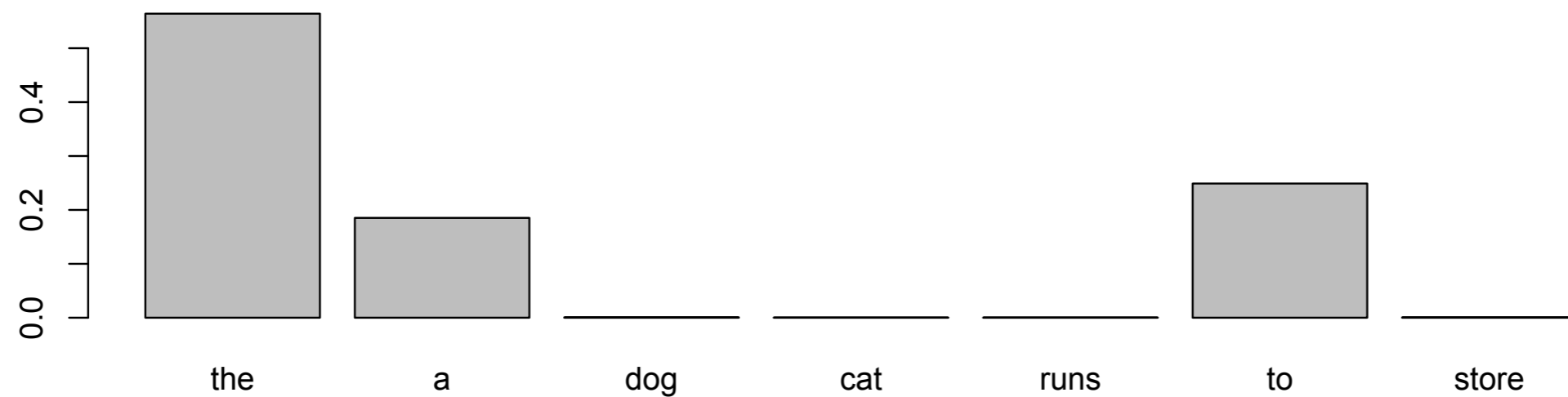
```
1 Given: a set  $\mathcal{X} = \{x_1, \dots, x_n\}$  of objects
2         a function  $\text{sim}: \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$ 
3 for  $i := 1$  to  $n$  do
4      $c_i := \{x_i\}$  end
5  $C := \{c_1, \dots, c_n\}$ 
6  $j := n + 1$ 
7 while  $C > 1$ 
8      $(c_{n_1}, c_{n_2}) := \arg \max_{(c_u, c_v) \in C \times C} \text{sim}(c_u, c_v)$ 
9      $c_j = c_{n_1} \cup c_{n_2}$ 
10     $C := C \setminus \{c_{n_1}, c_{n_2}\} \cup \{c_j\}$ 
11     $j := j + 1$ 
```

Similarity

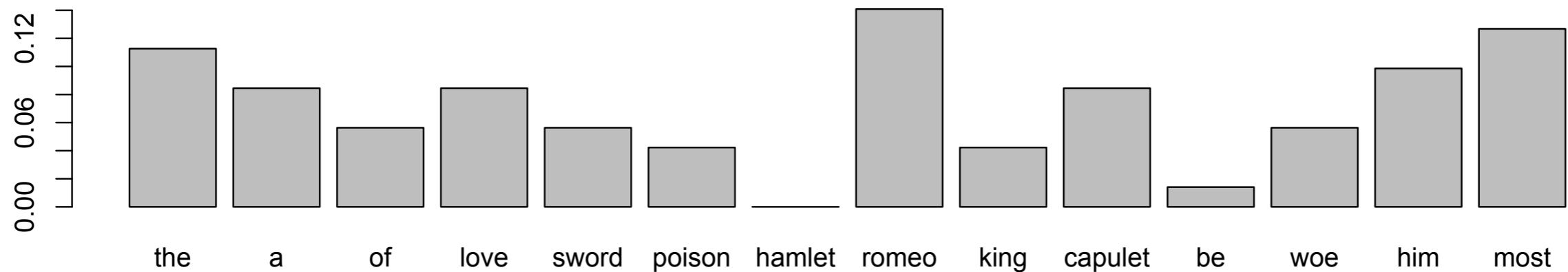
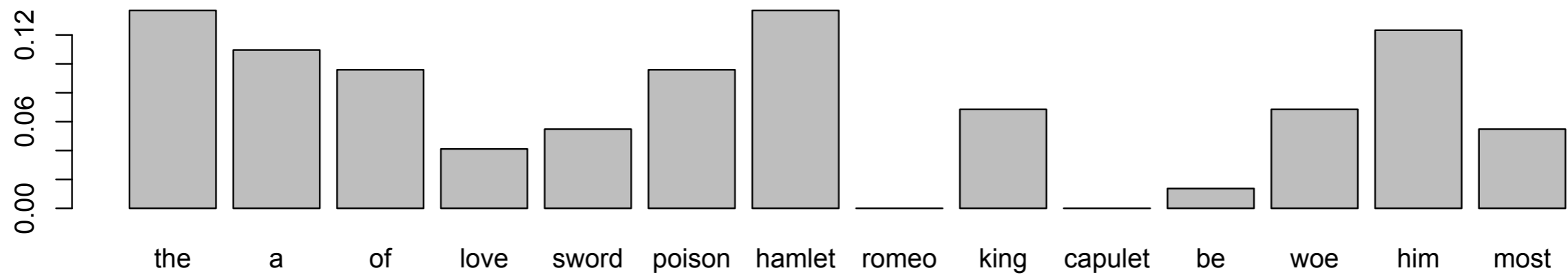
$$\mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{X}) \rightarrow \mathbb{R}$$

- What are you comparing?
- How do you quantify the similarity/difference of those things?

Probability



Unigram probability



Similarity

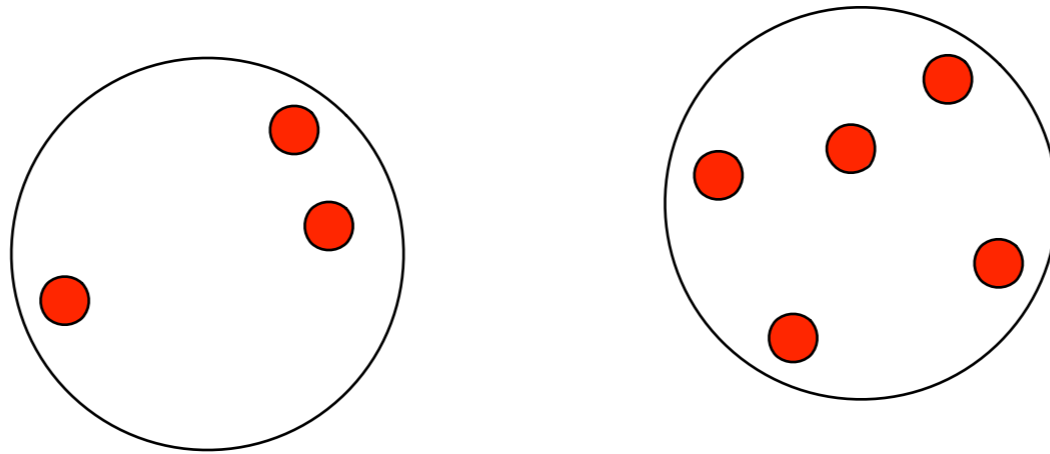
$$\text{Euclidean} = \sqrt{\sum_i^{vocab} (P_i^{\text{Hamlet}} - P_i^{\text{Romeo}})^2}$$

Cosine similarity, Jensen-Shannon divergence...

Cluster similarity



Cluster similarity



- Single link: two **most** similar elements
- Complete link: two **least** similar elements
- Group average: average of all members

Flat Clustering

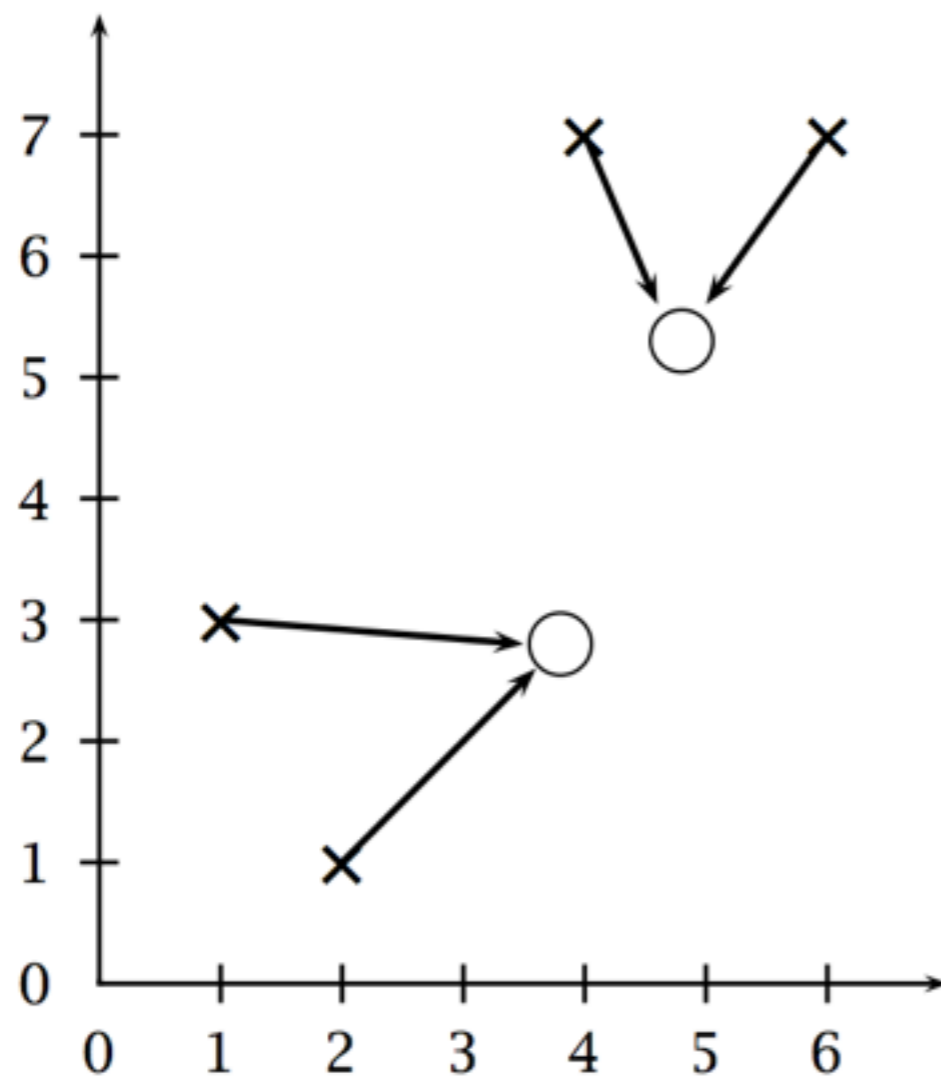
- Partitions the data into a set of K clusters



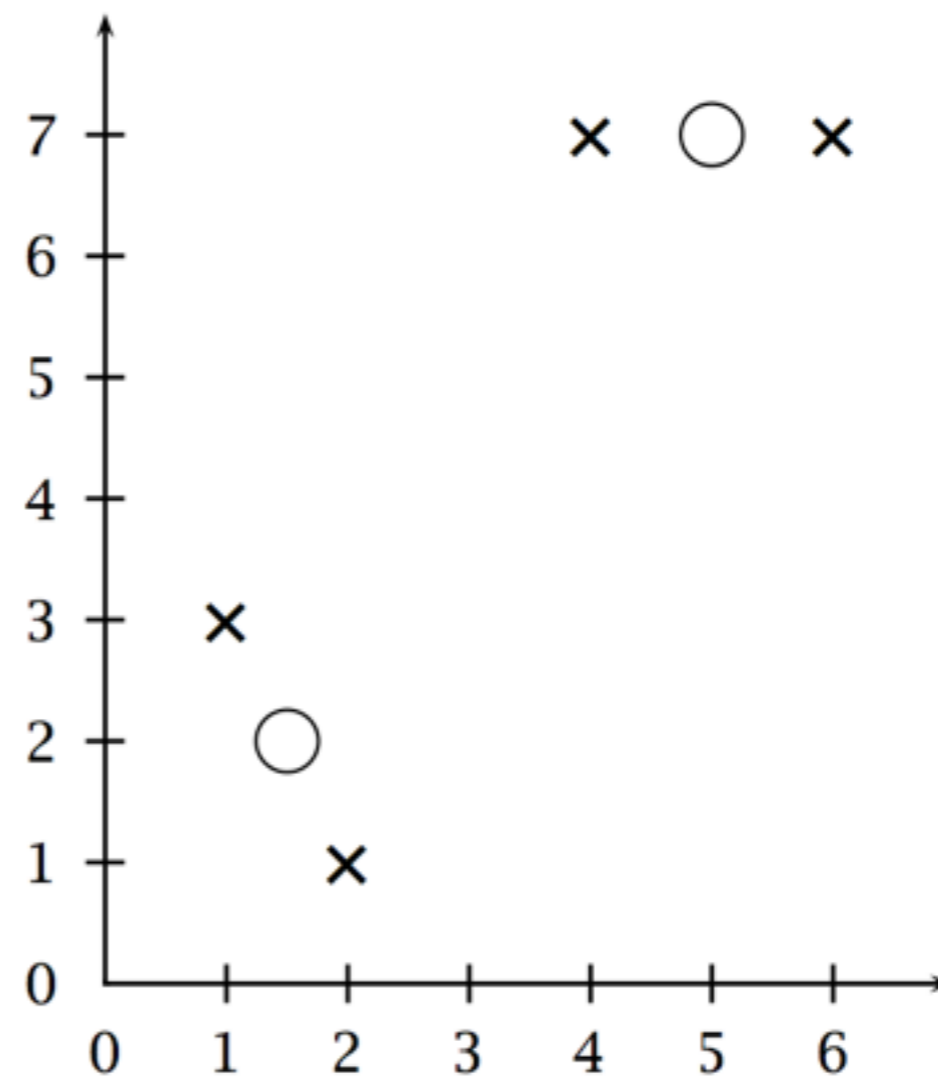
K-means

```
1 Given: a set  $\mathcal{X} = \{\vec{x}_1, \dots, \vec{x}_n\} \subseteq \mathbb{R}^m$ 
2         a distance measure  $d : \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}$ 
3         a function for computing the mean  $\mu : \mathcal{P}(\mathbb{R}) \rightarrow \mathbb{R}^m$ 
4 Select  $k$  initial centers  $\vec{f}_1, \dots, \vec{f}_k$ 
5 while stopping criterion is not true do
6     for all clusters  $c_j$  do
7          $c_j = \{\vec{x}_i \mid \forall \vec{f}_l \ d(\vec{x}_i, \vec{f}_j) \leq d(\vec{x}_i, \vec{f}_l)\}$ 
8     end
9     for all means  $\vec{f}_j$  do
10          $\vec{f}_j = \mu(c_j)$ 
11     end
12 end
```

K-means



assignment

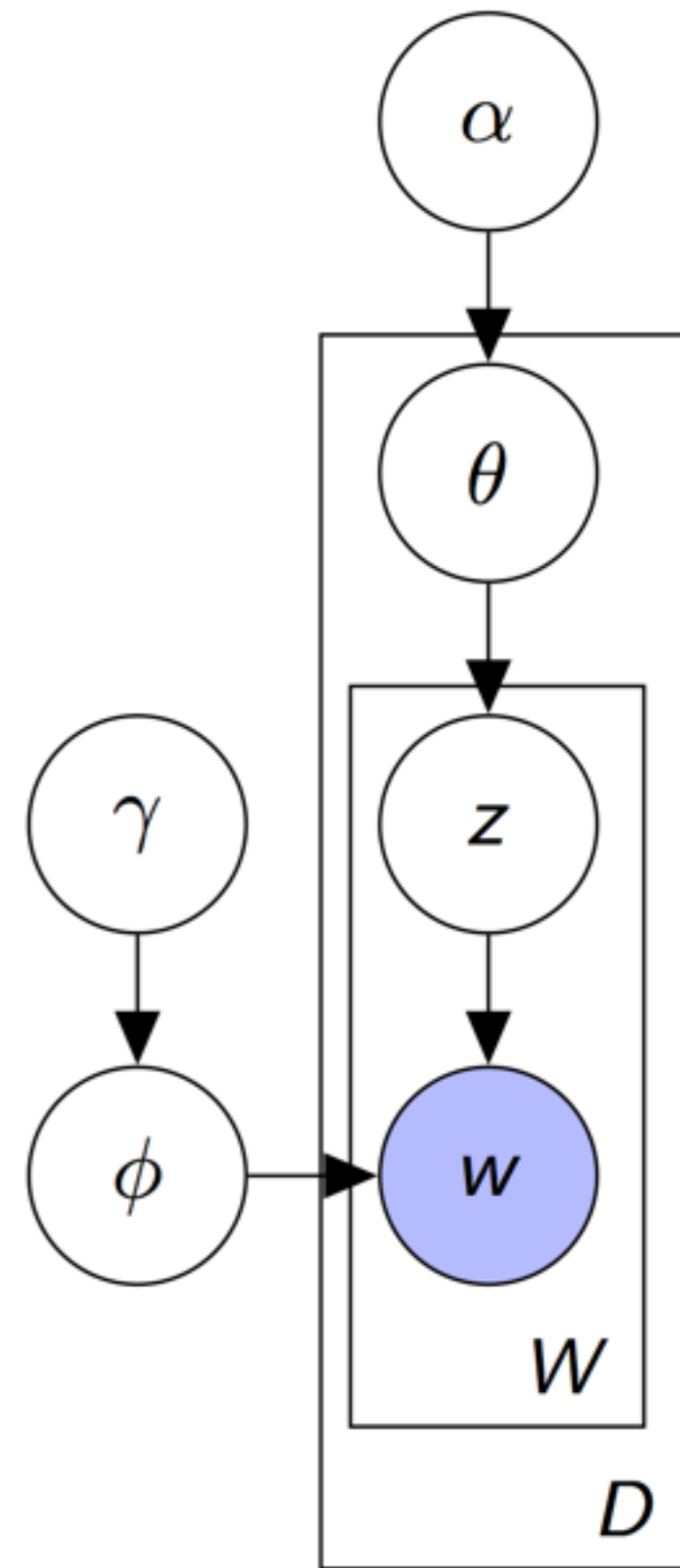


recomputation of means

Try it yourself

- Shakespeare + English stoplist
<http://bit.ly/1hdKX0R>
- <http://lexos.wheatoncollege.edu>

Topic Models



Topic Models

- A probabilistic model for discovering hidden “topics” or “themes” (groups of terms that tend to occur together) in documents.
- Unsupervised (find *interesting structure* in the data)
- Clustering algorithm

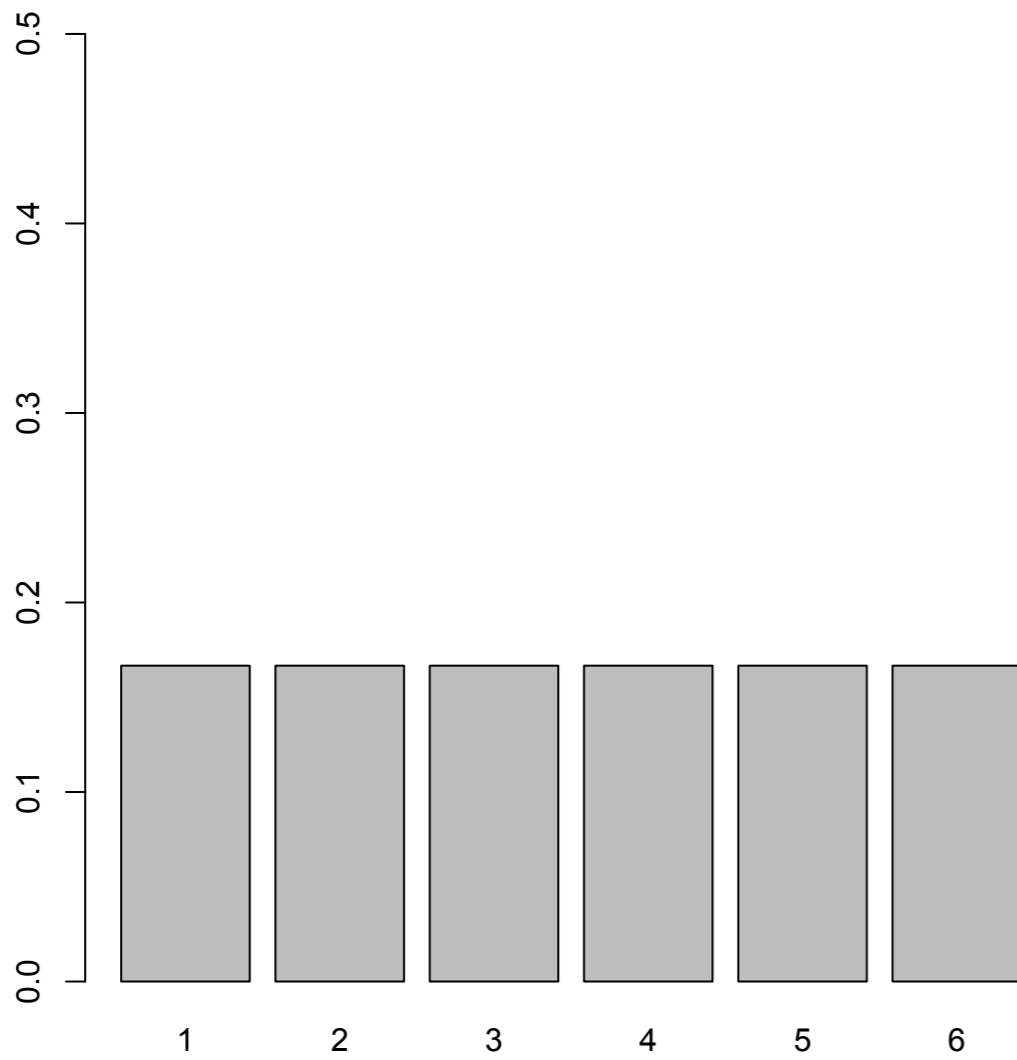
Topic Models

- **Input:** set of documents, number of clusters to learn.
- **Output:**
 - topics
 - topic ratio in each document
 - topic distribution for each word in doc

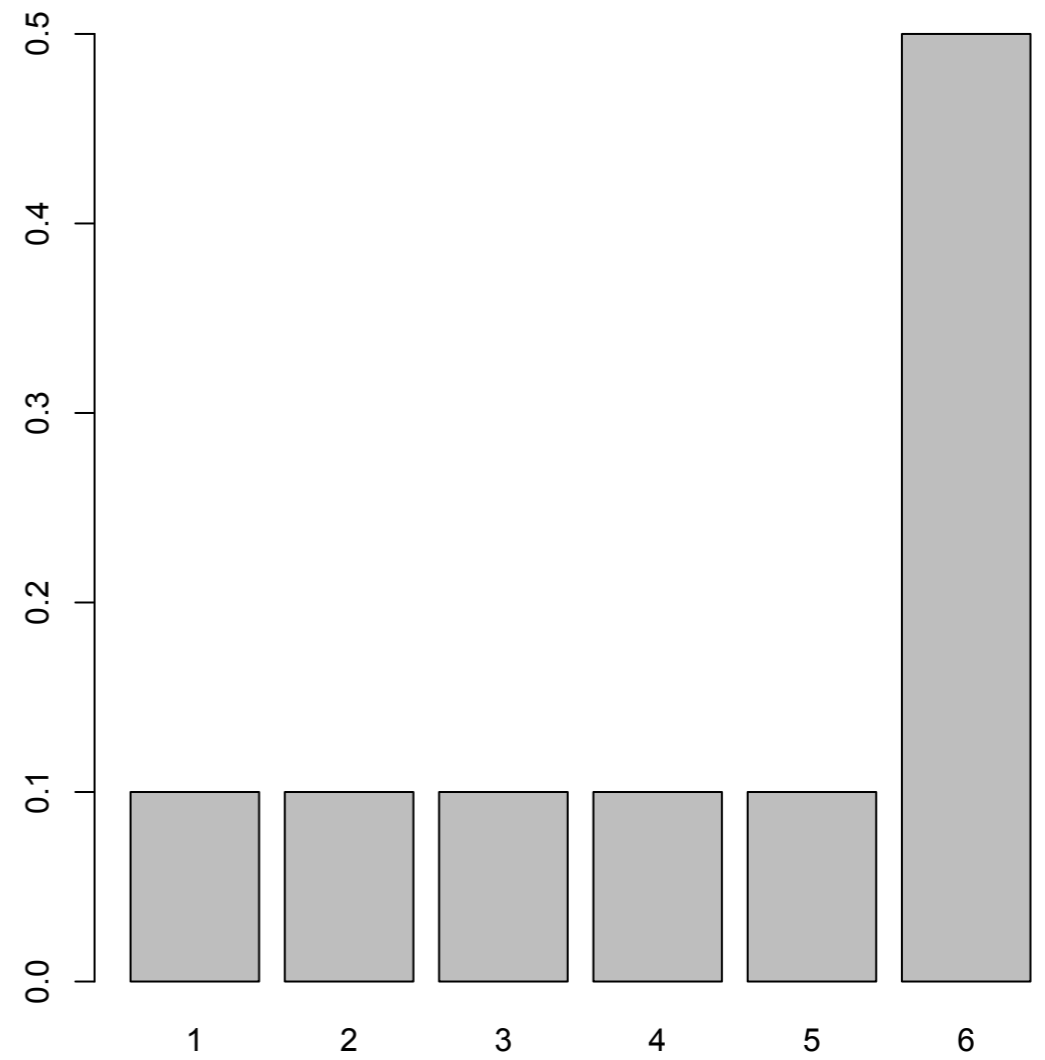
{album, band, music}	{government, party, election}	{game, team, player}
album	government	game
band	party	team
music	election	player
song	state	win
release	political	play
{god, call, give}	{company, market, business}	{math, number, function}
god	company	math
call	market	number
give	business	function
man	year	code
time	product	set
{city, large, area}	{math, energy, light}	{law, state, case}
city	math	law
large	energy	state
area	light	case
station	field	court
include	star	legal

Probability

fair

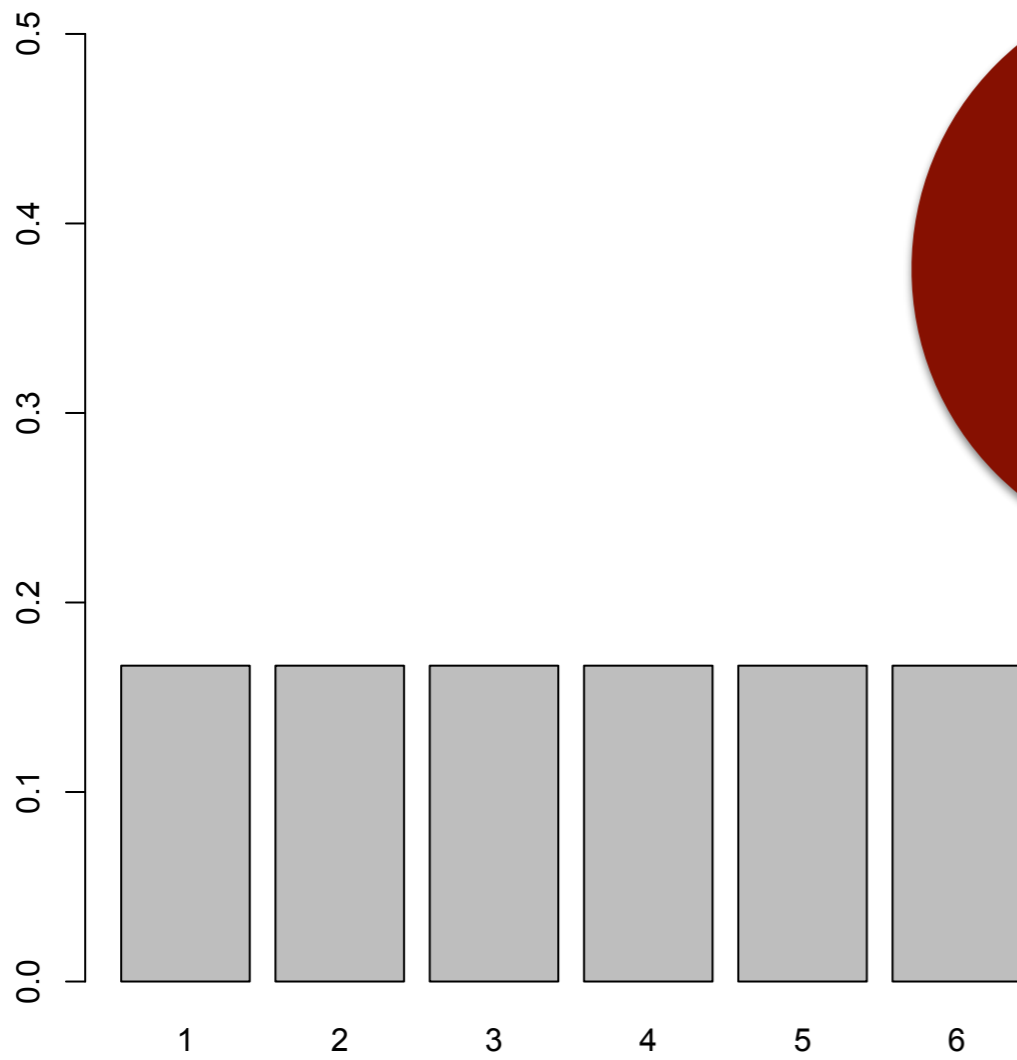


not fair

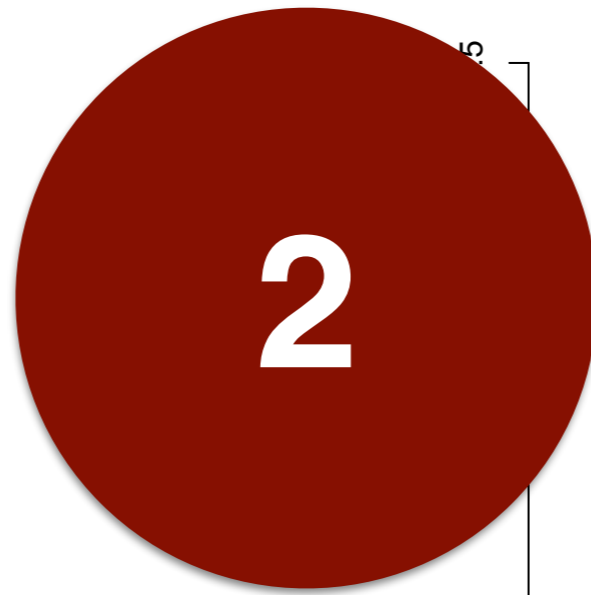
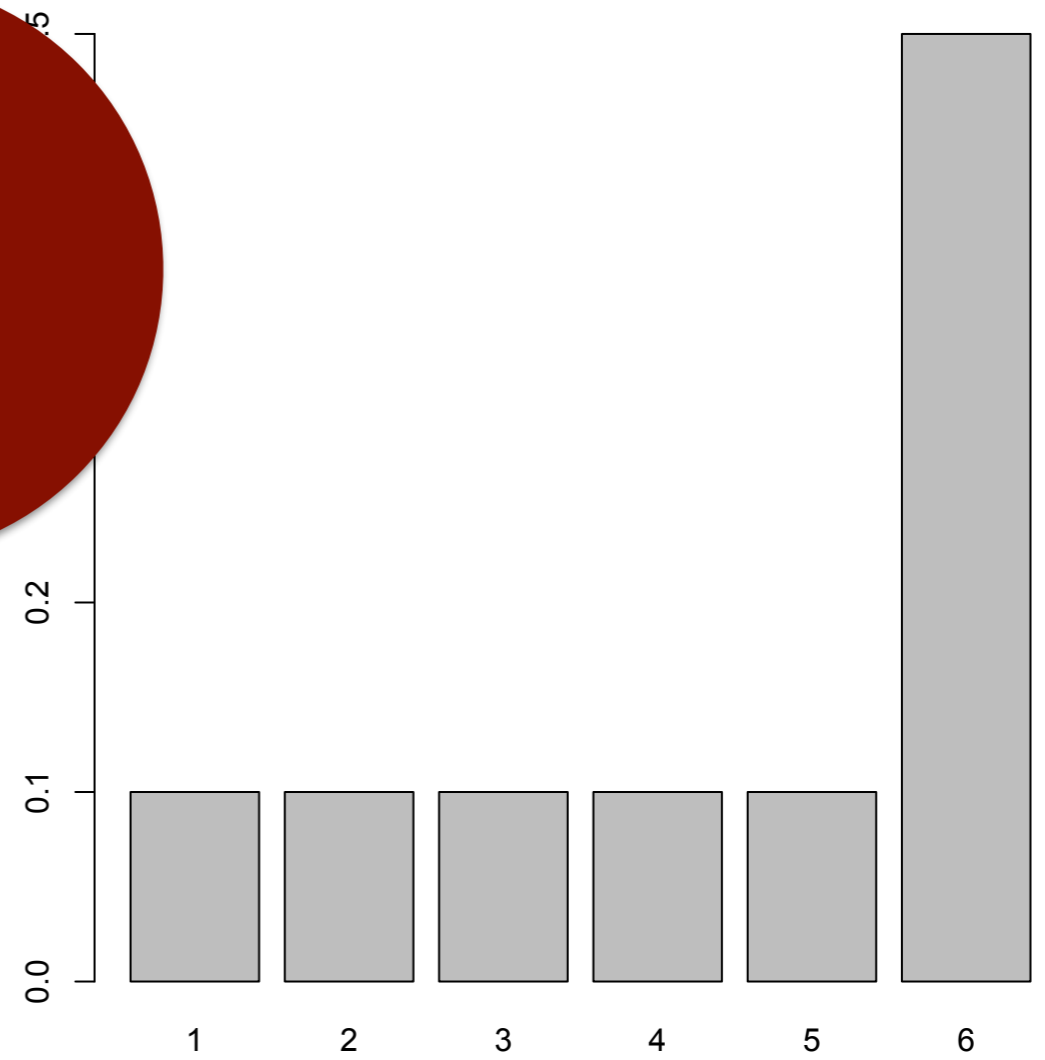


Probability

fair



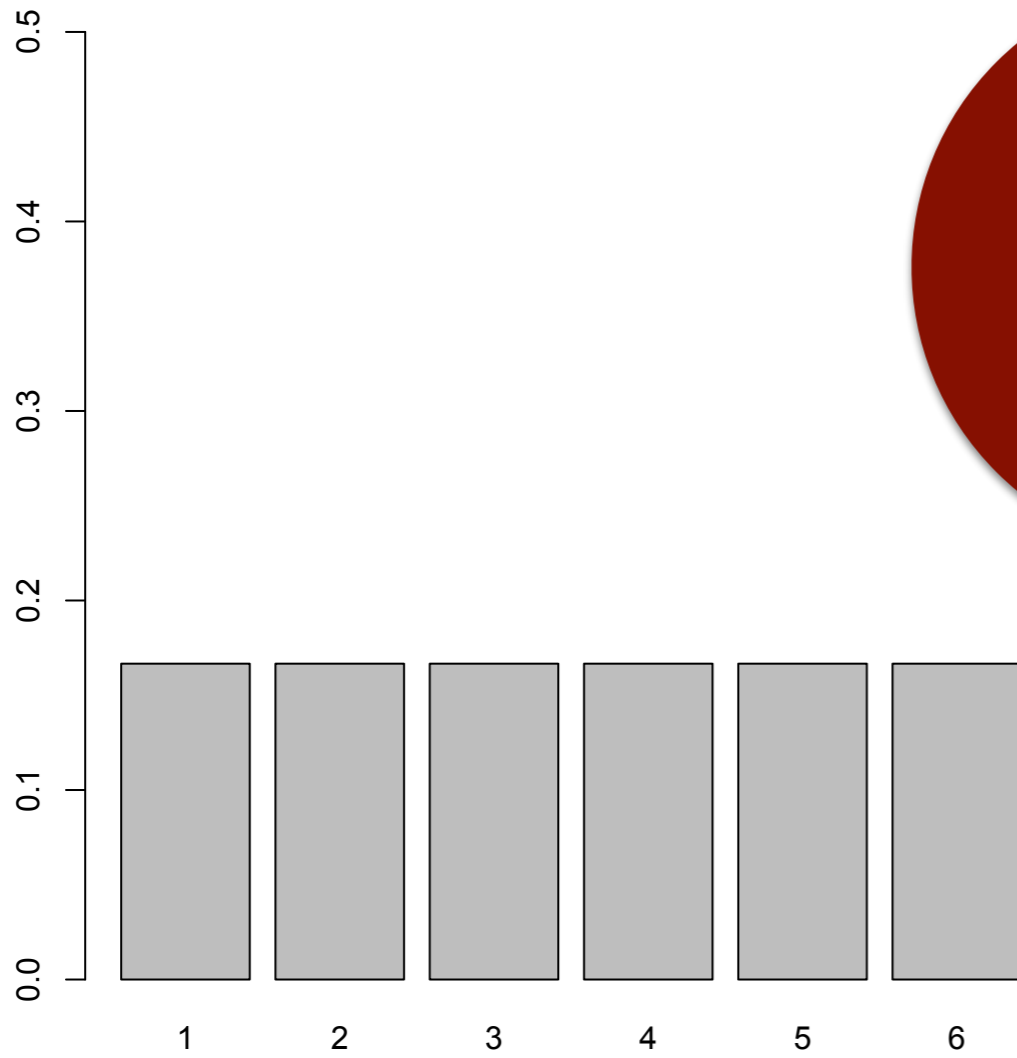
not fair



Probability

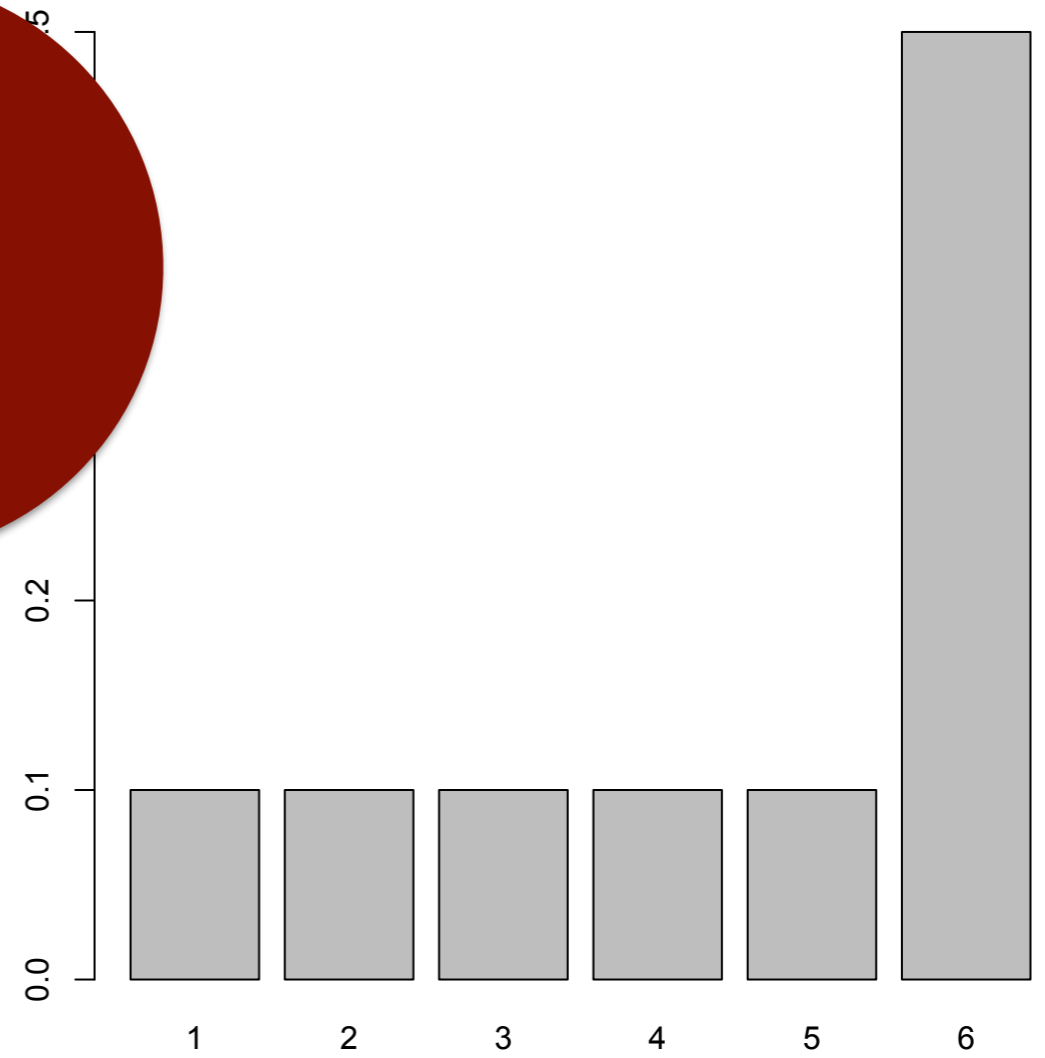
2

fair

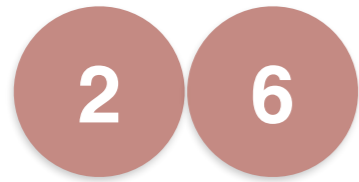


6

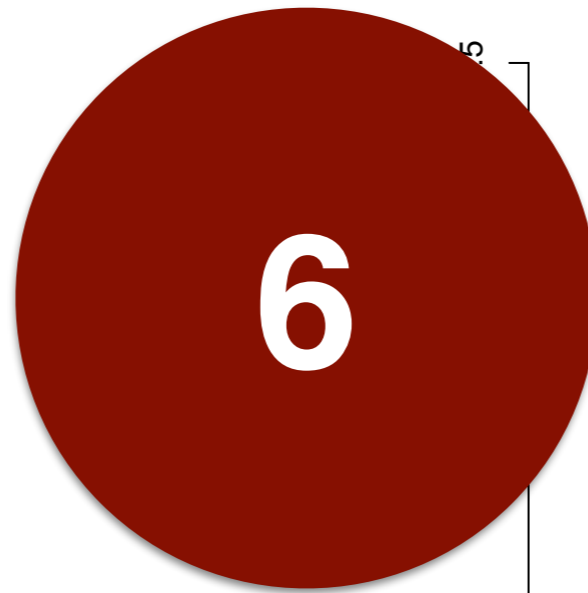
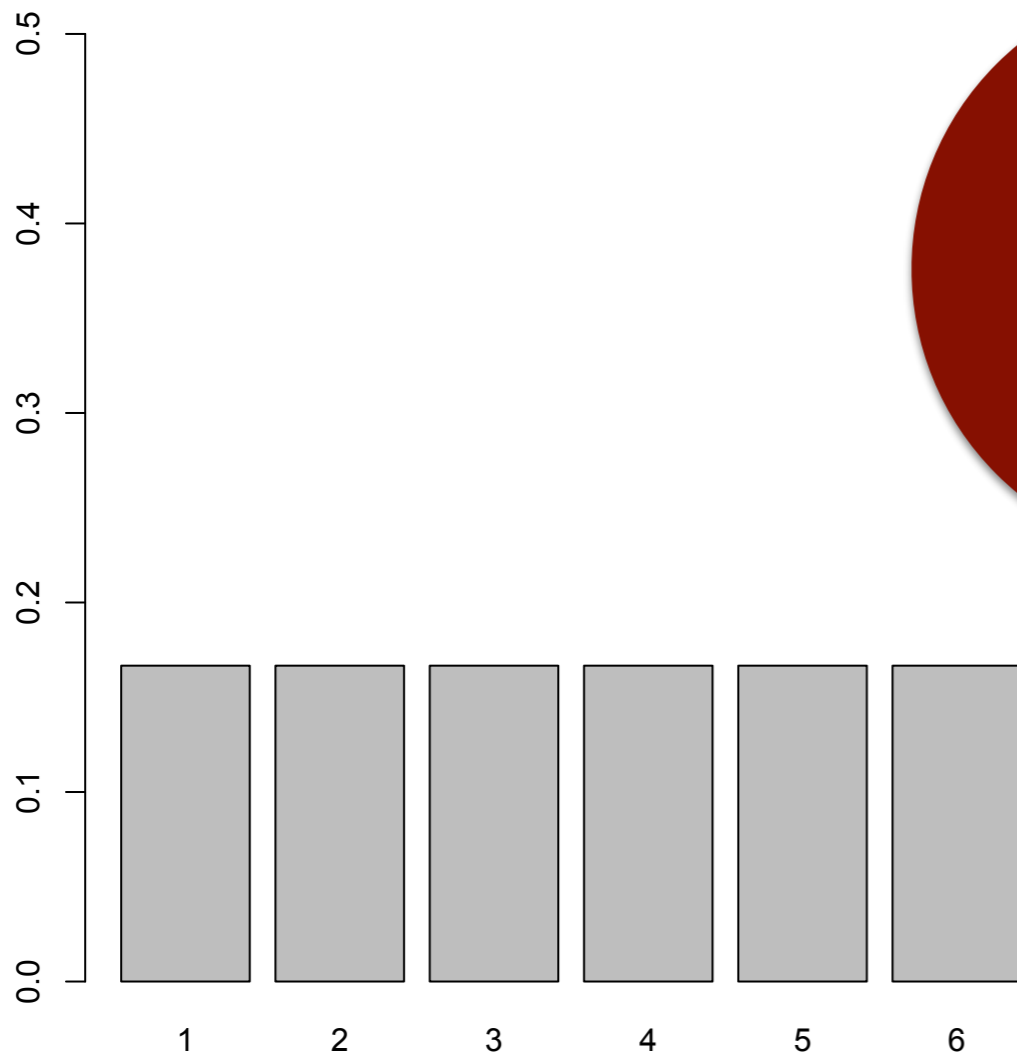
not fair



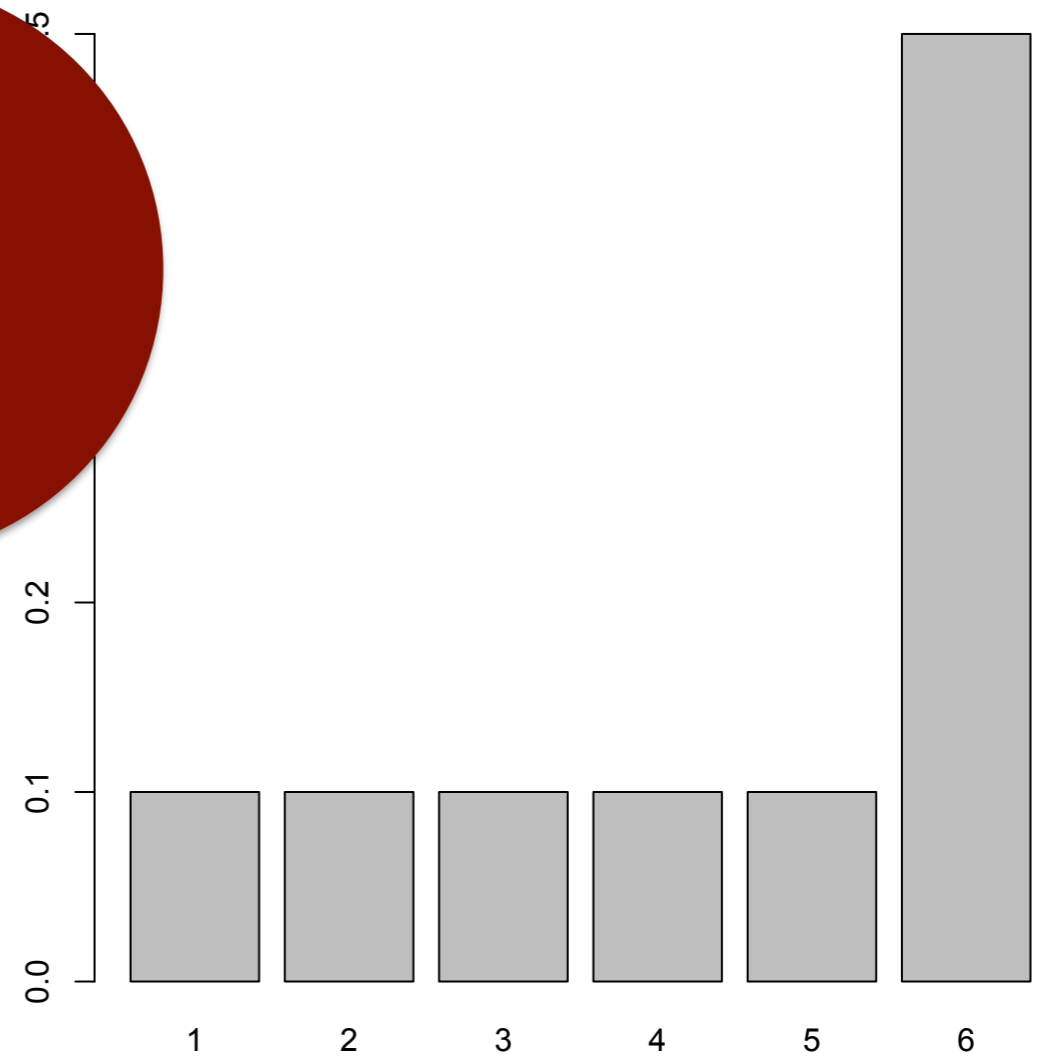
Probability



fair



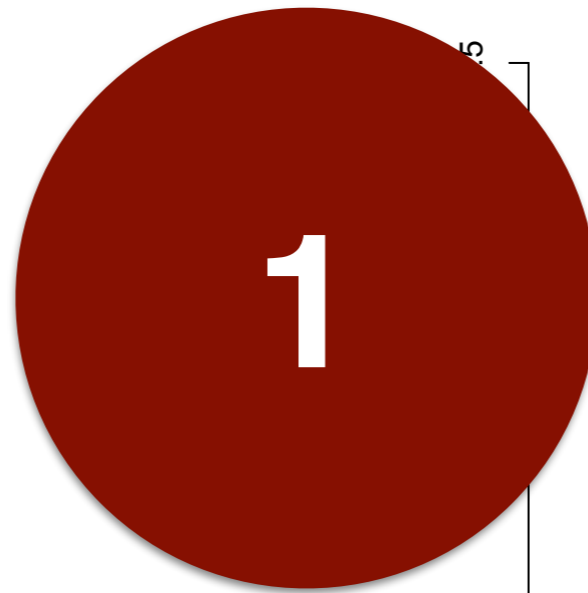
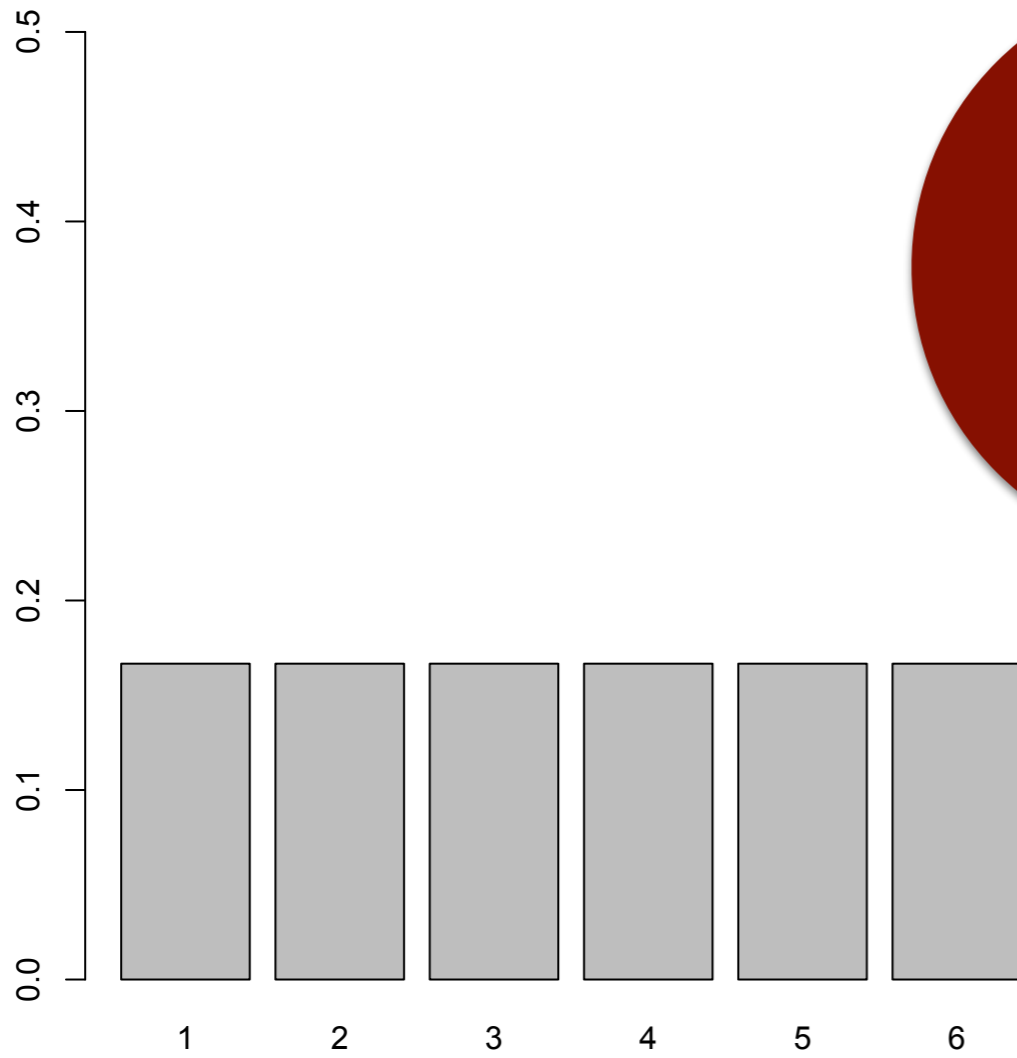
not fair



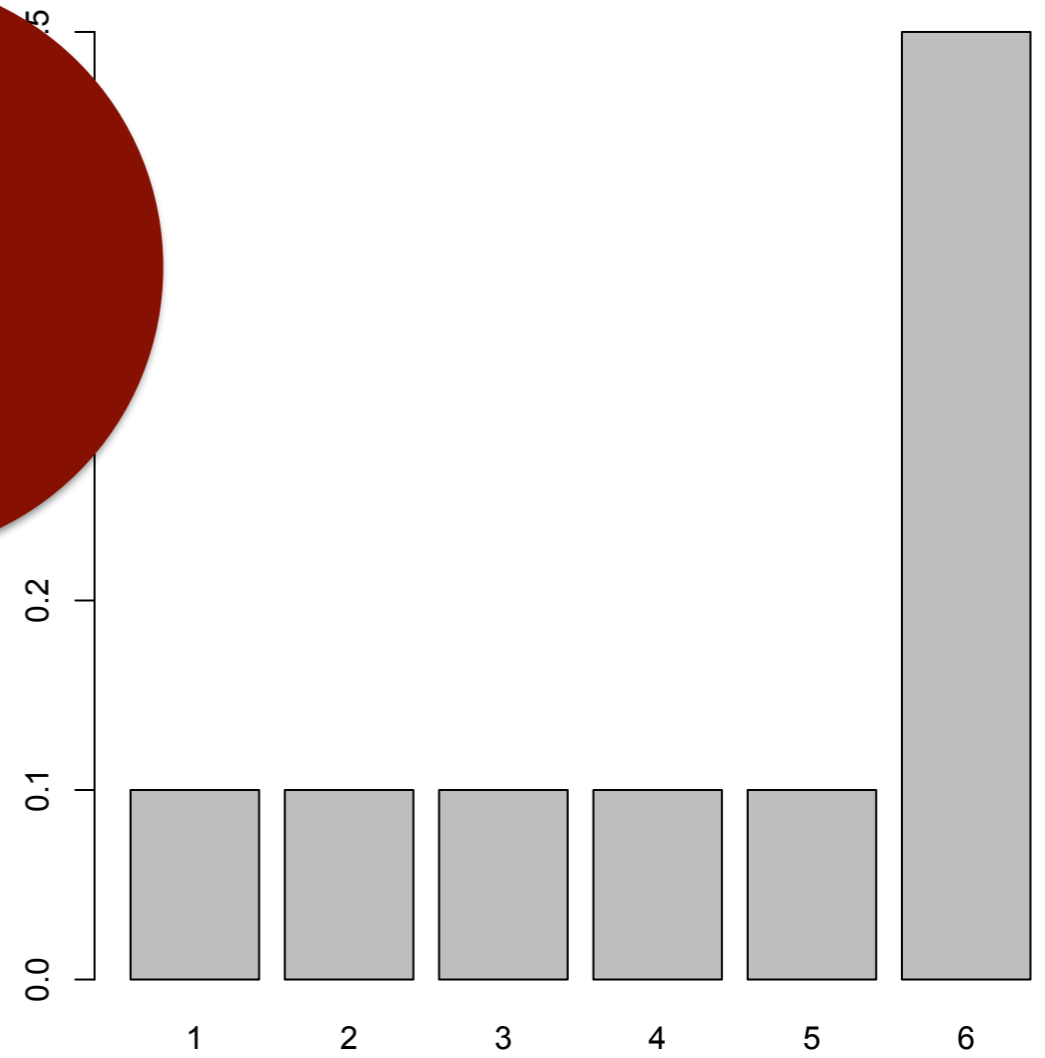
Probability



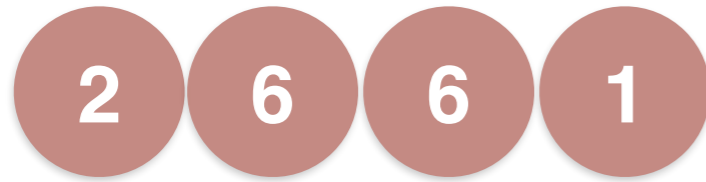
fair



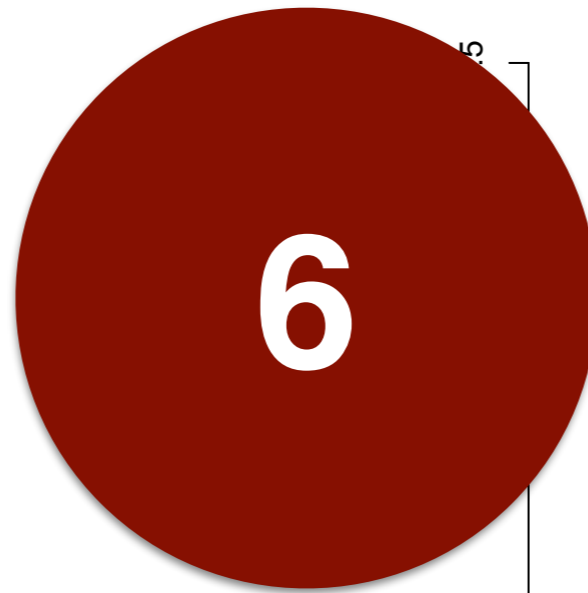
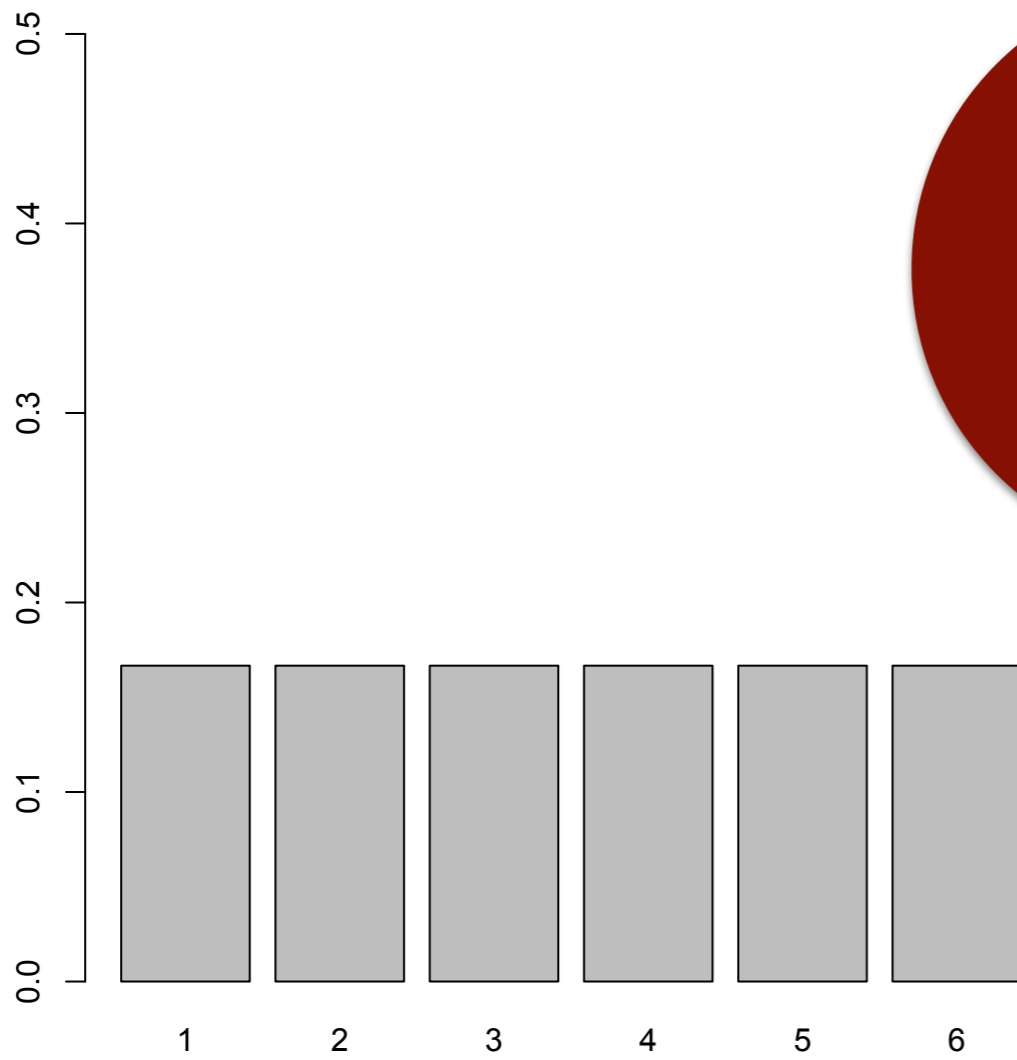
not fair



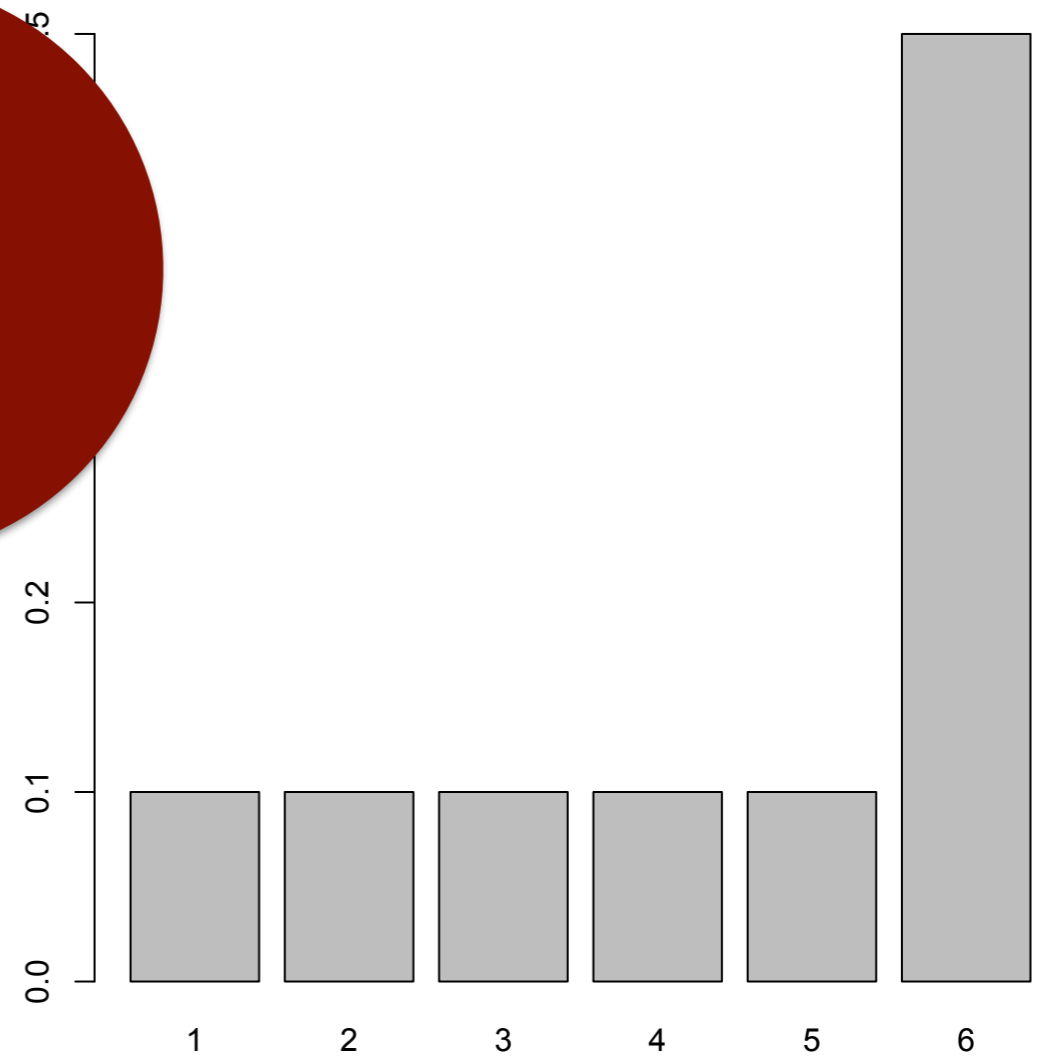
Probability



fair



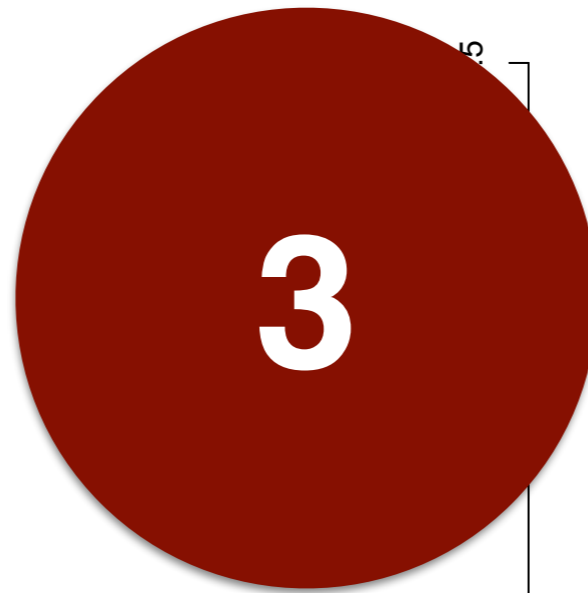
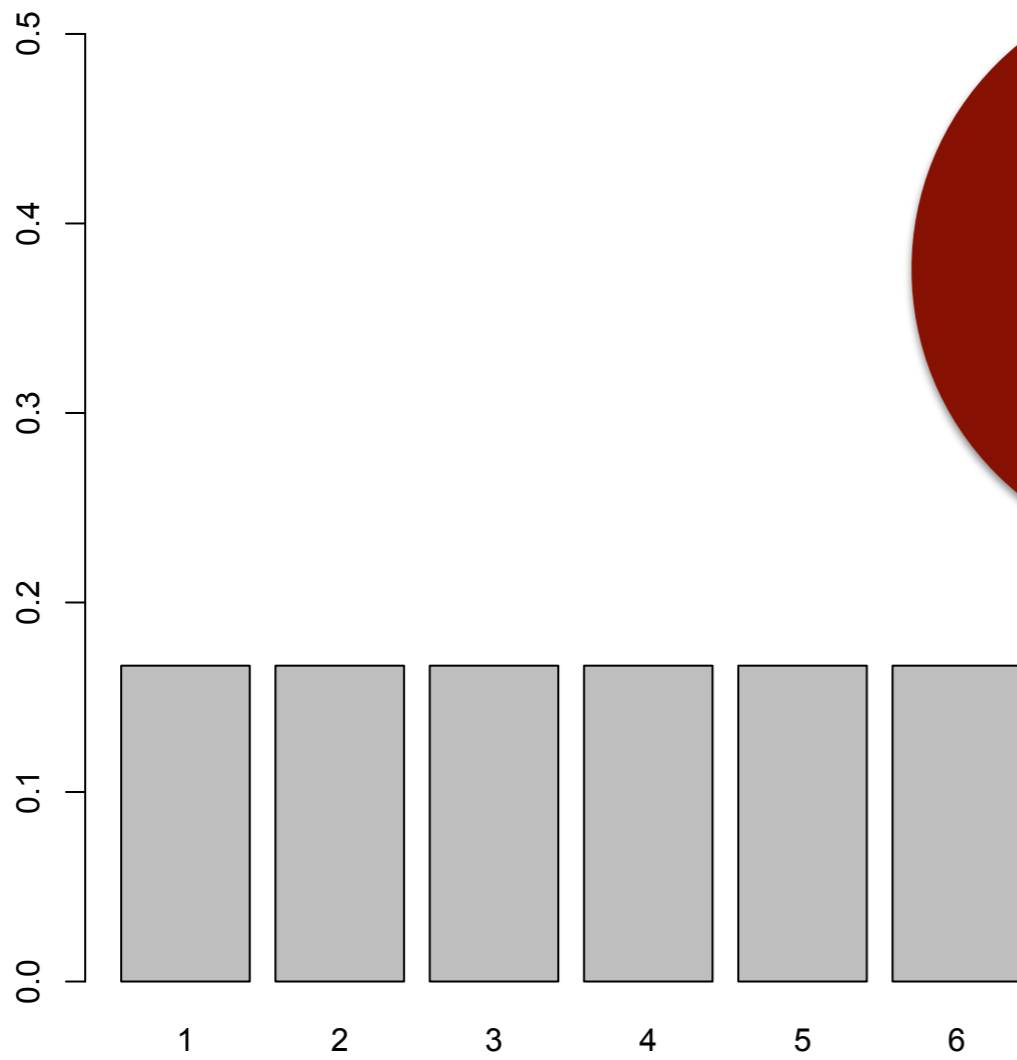
not fair



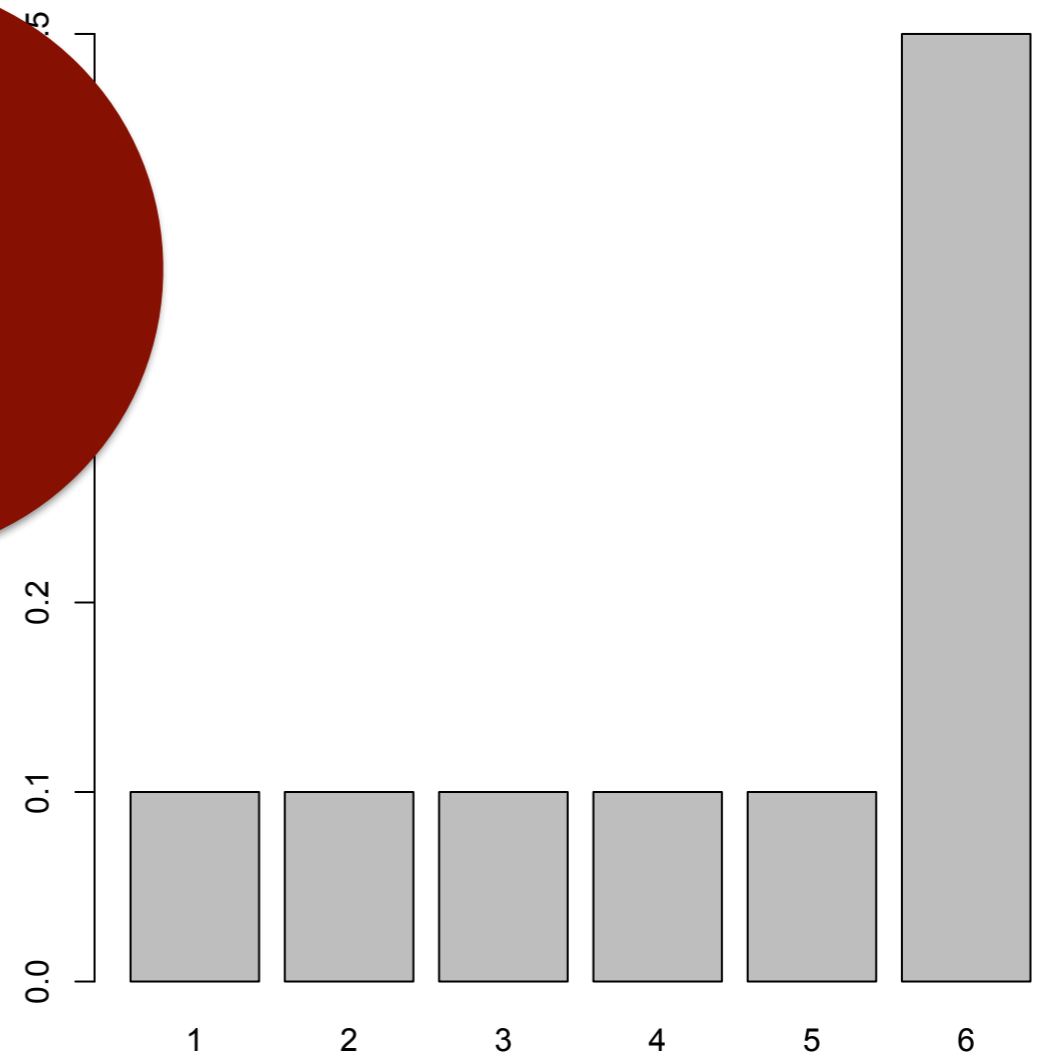
Probability



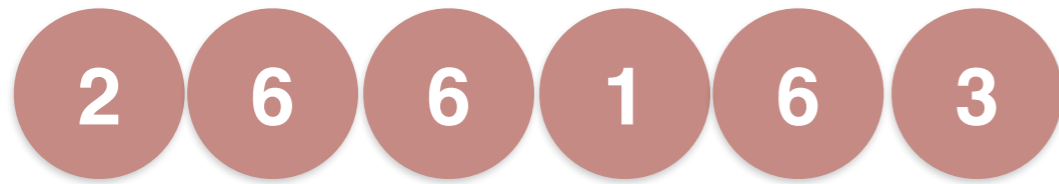
fair



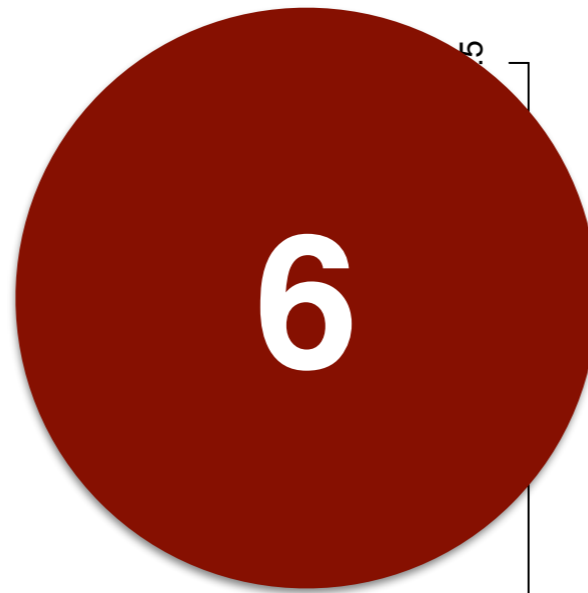
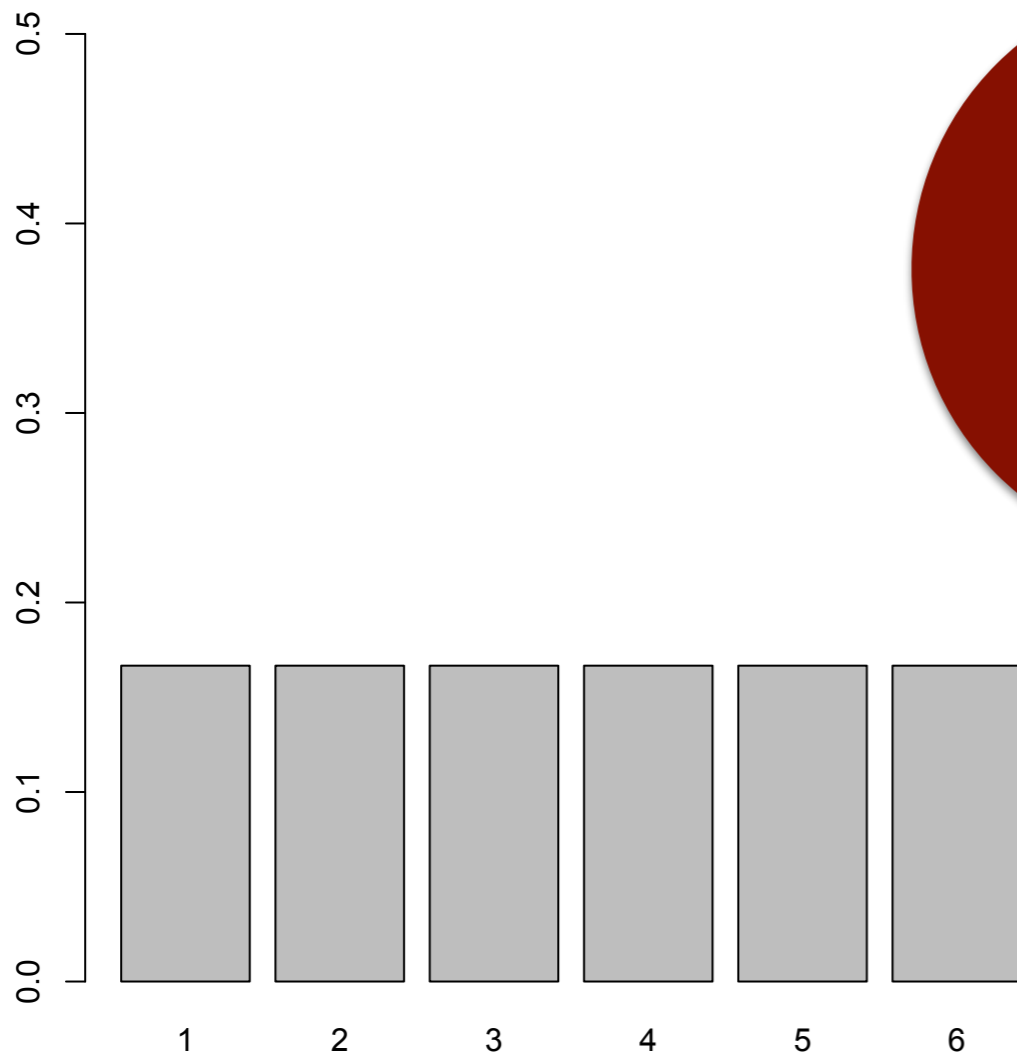
not fair



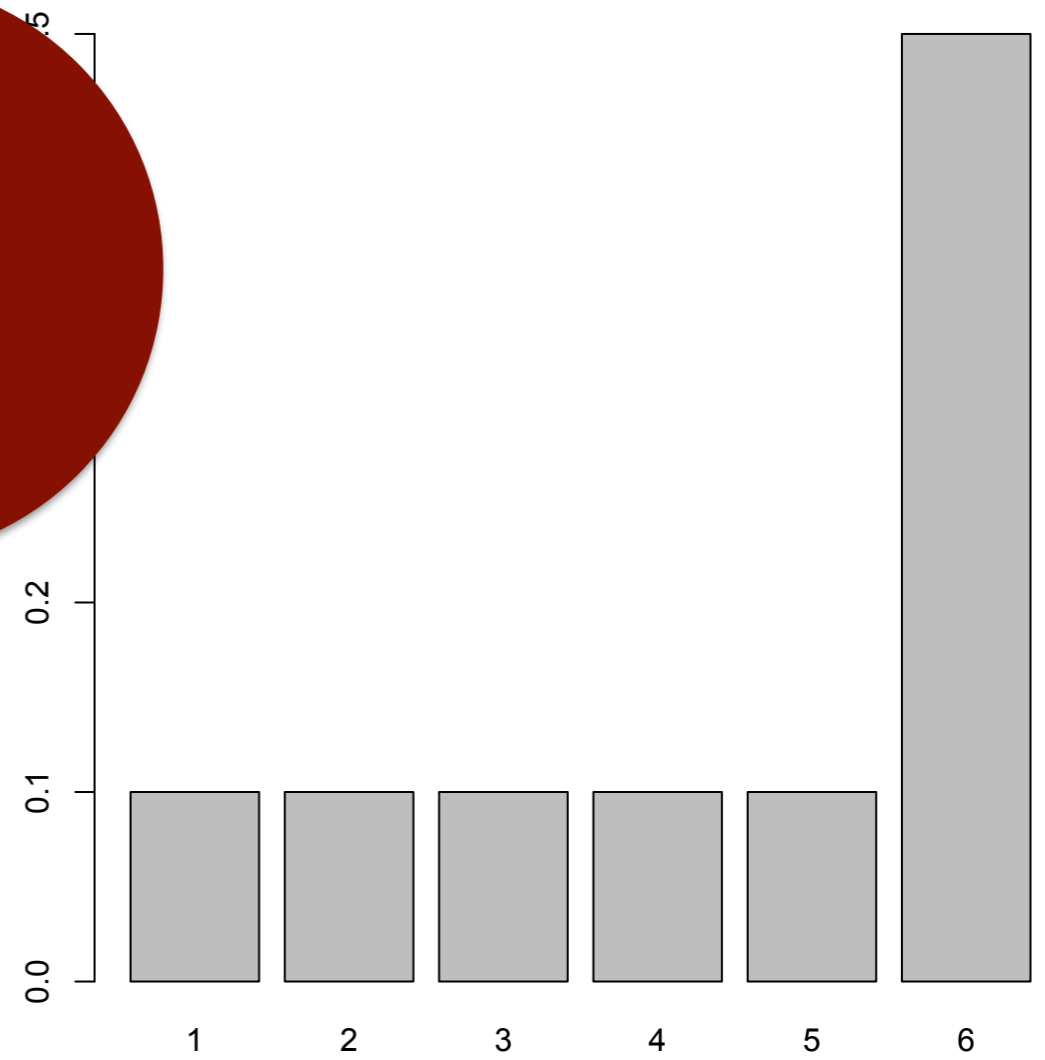
Probability



fair



not fair

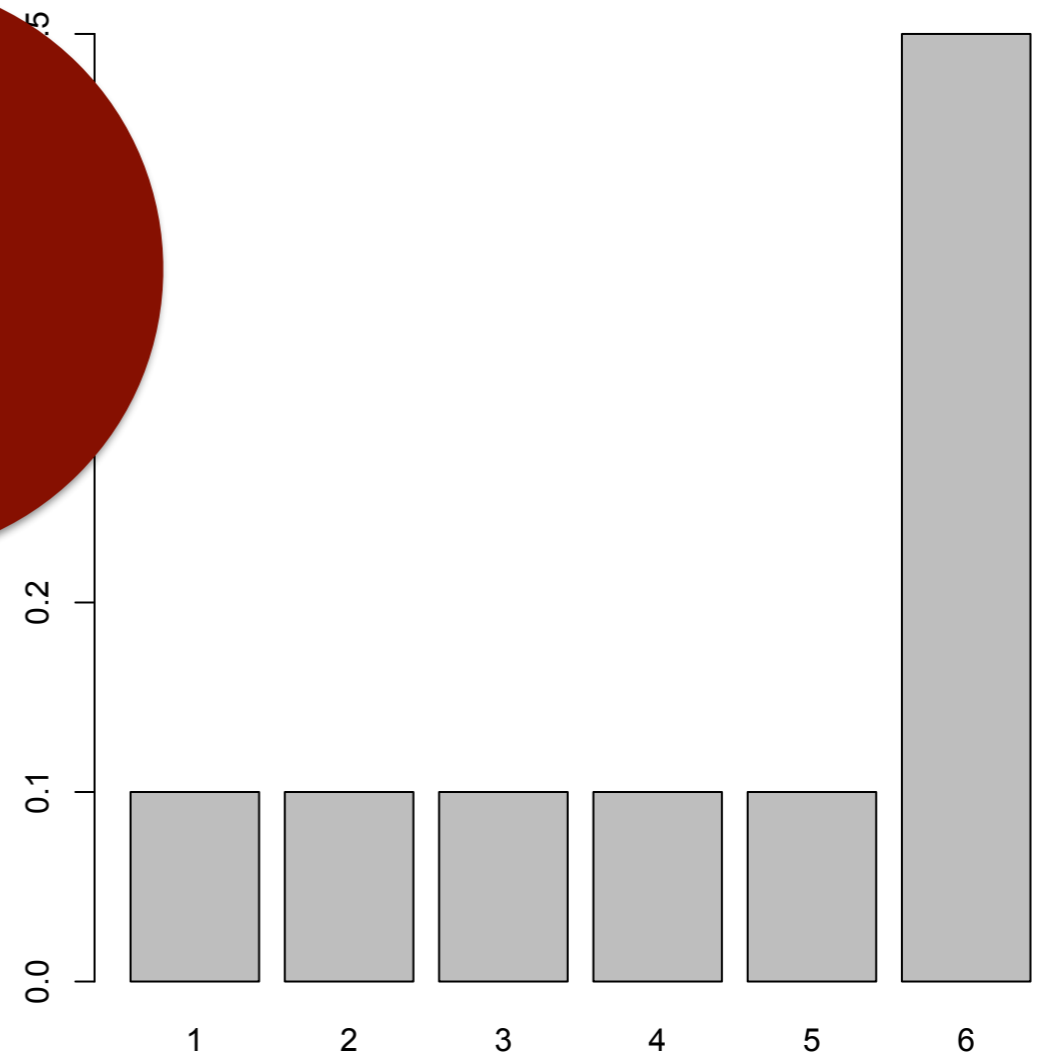
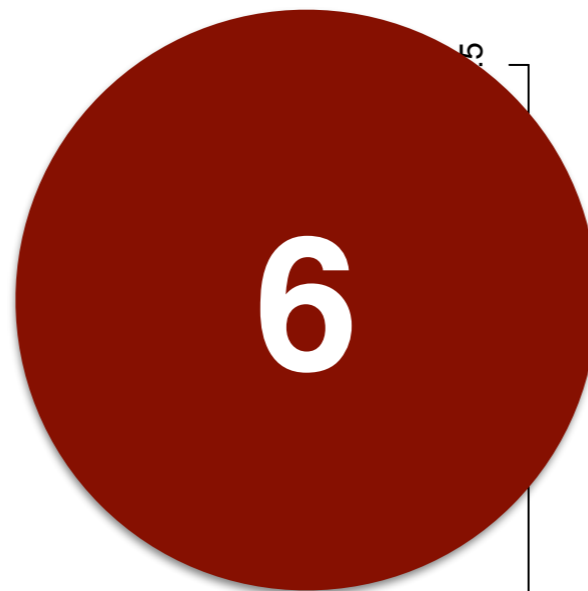
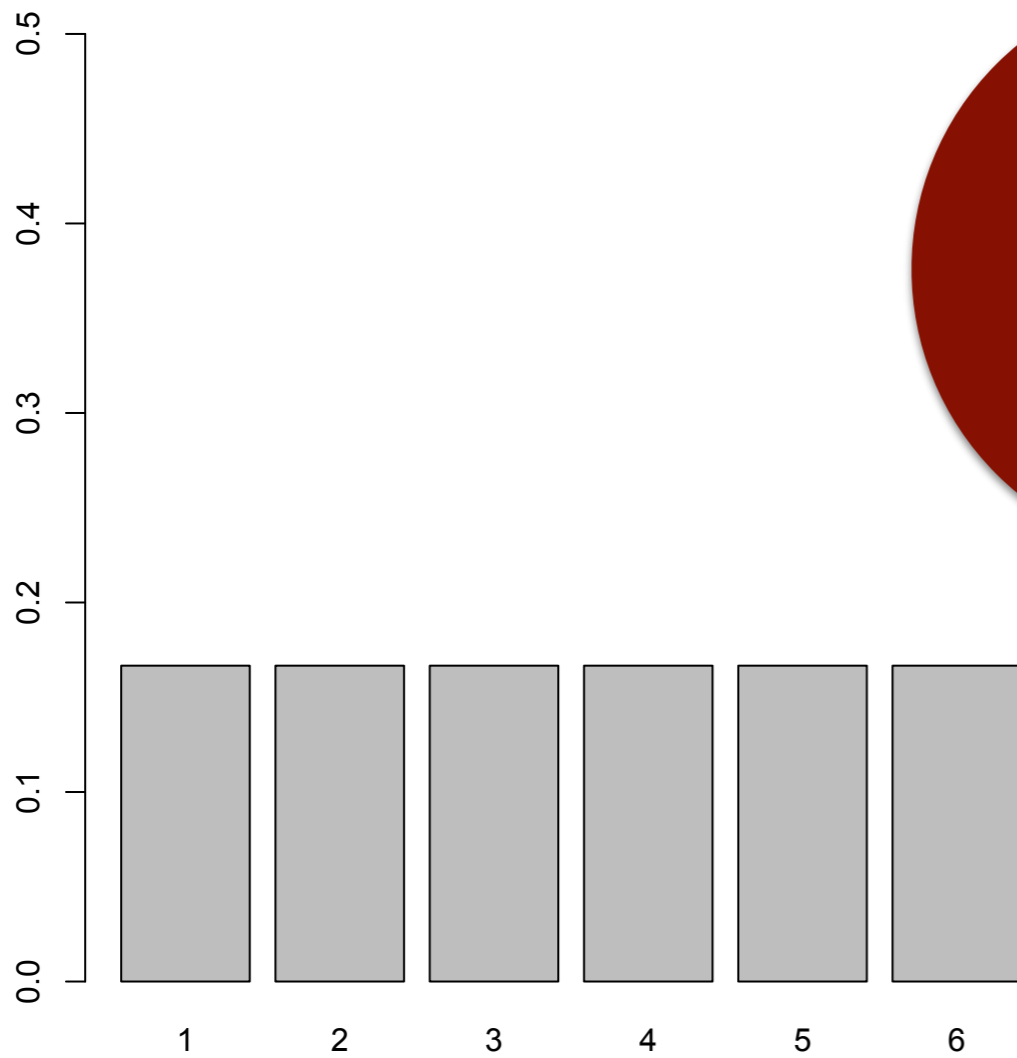


Probability

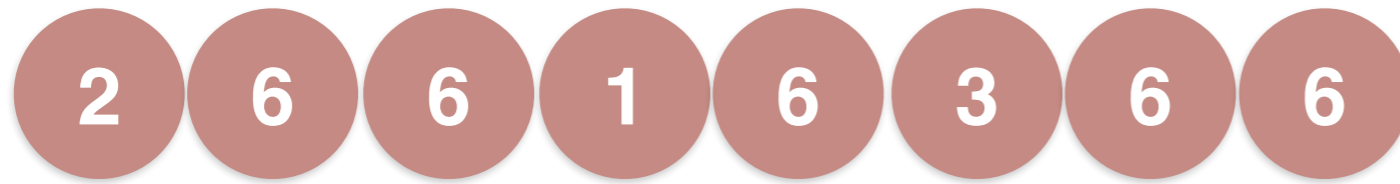


fair

not fair

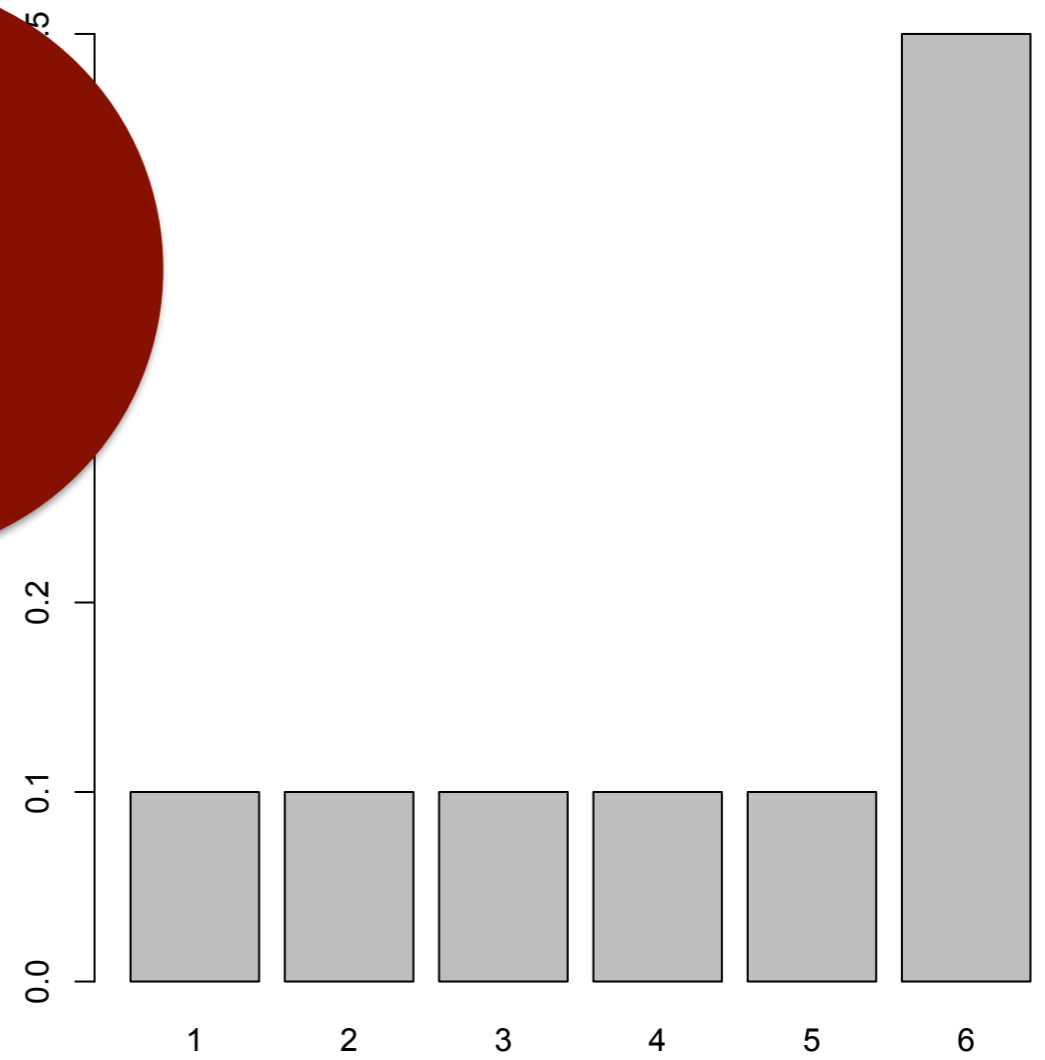
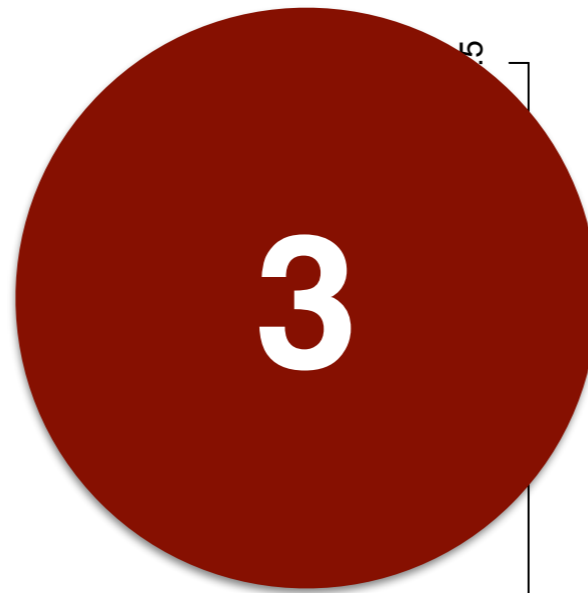
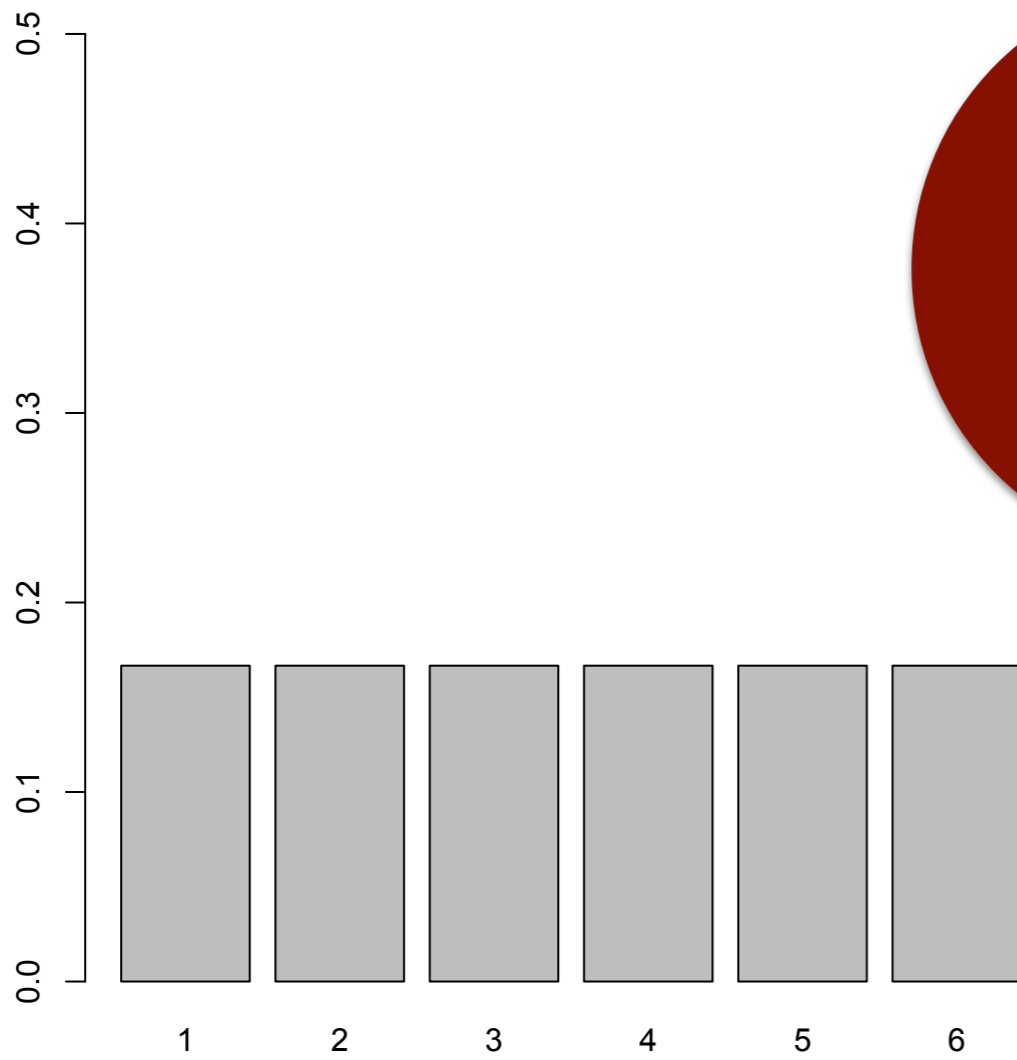


Probability



fair

not fair

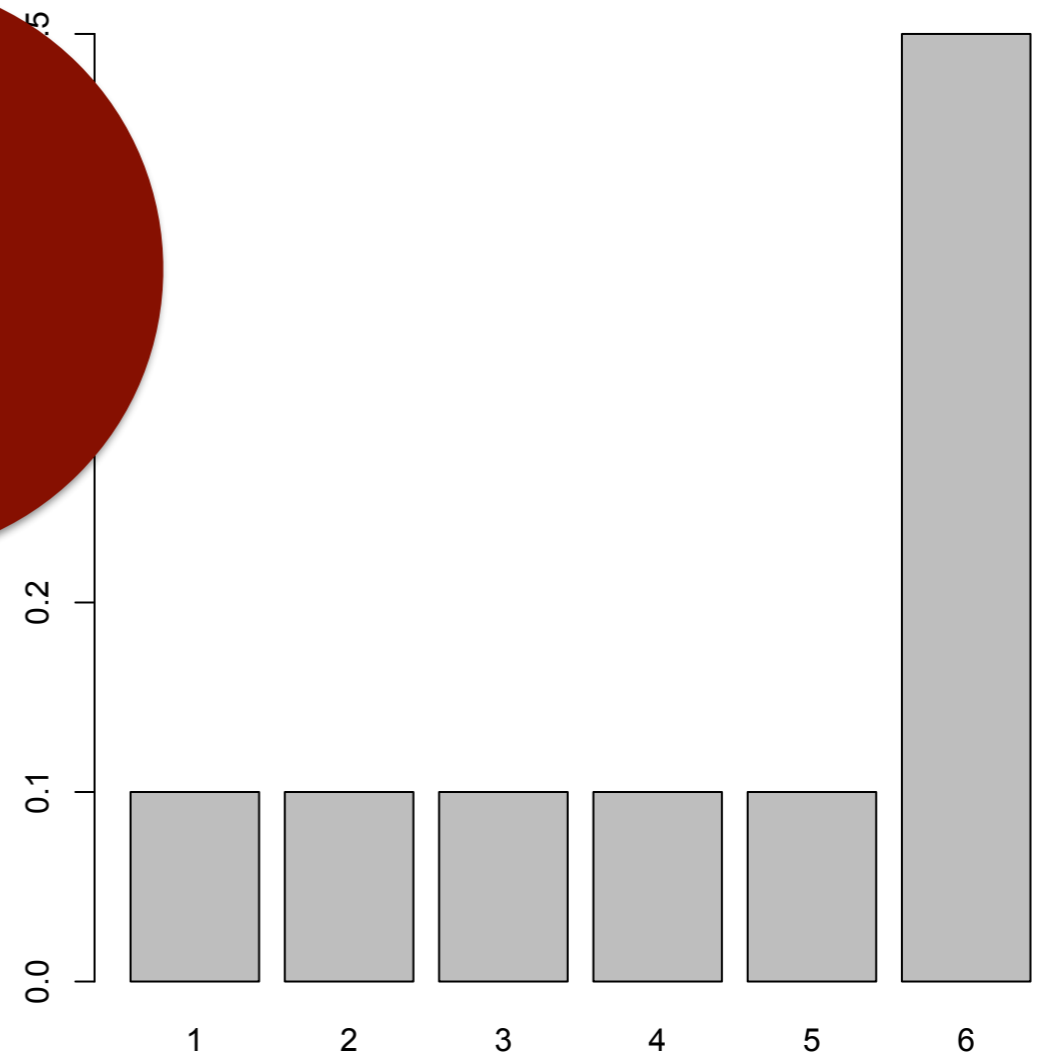
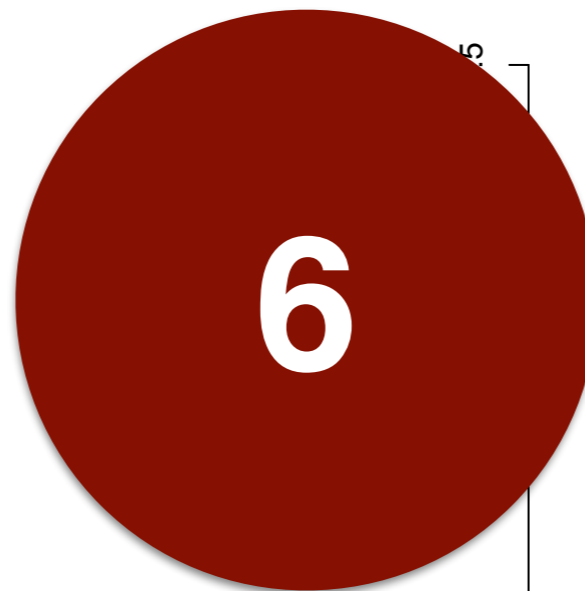
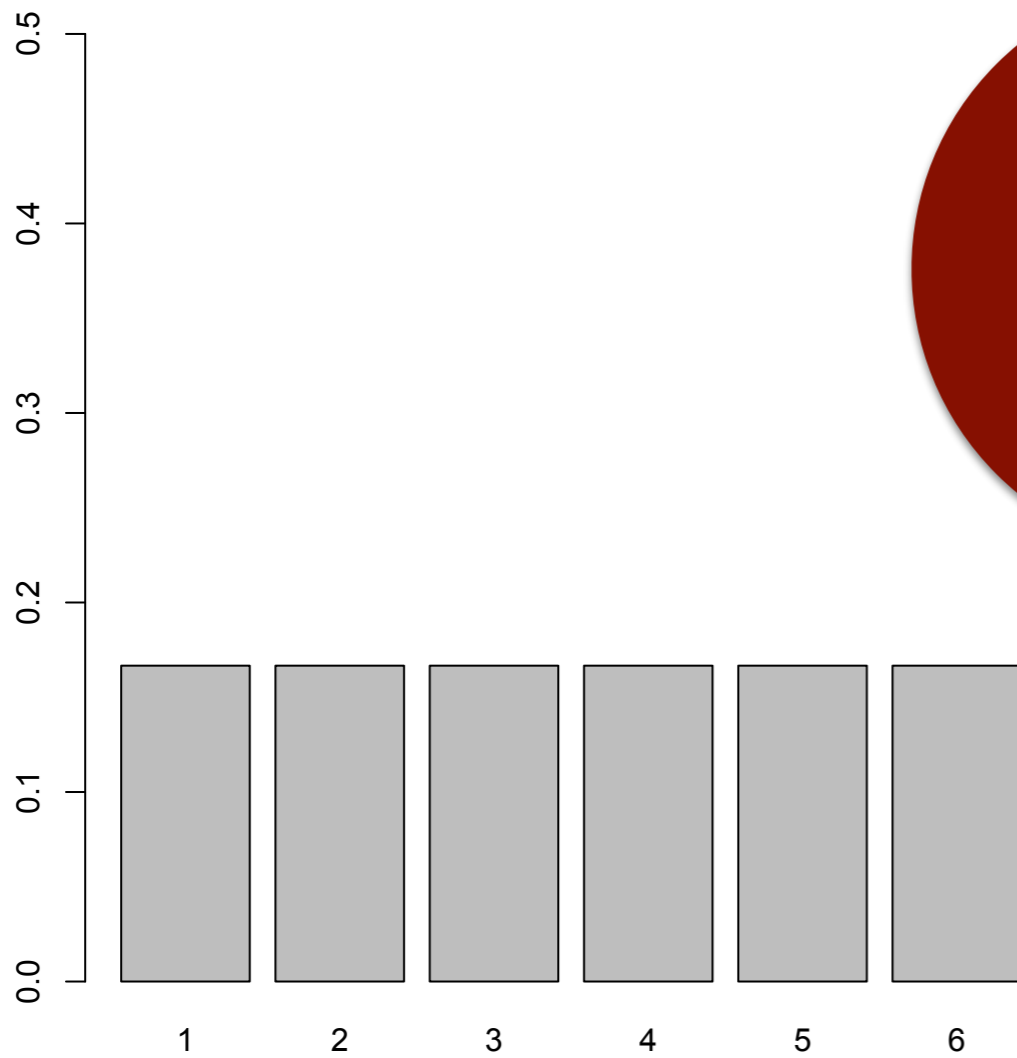


Probability

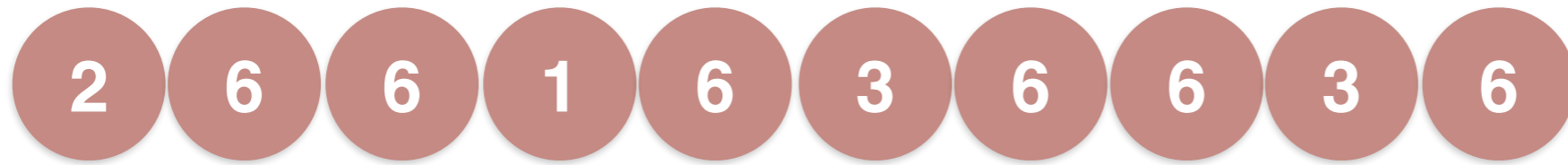


fair

not fair

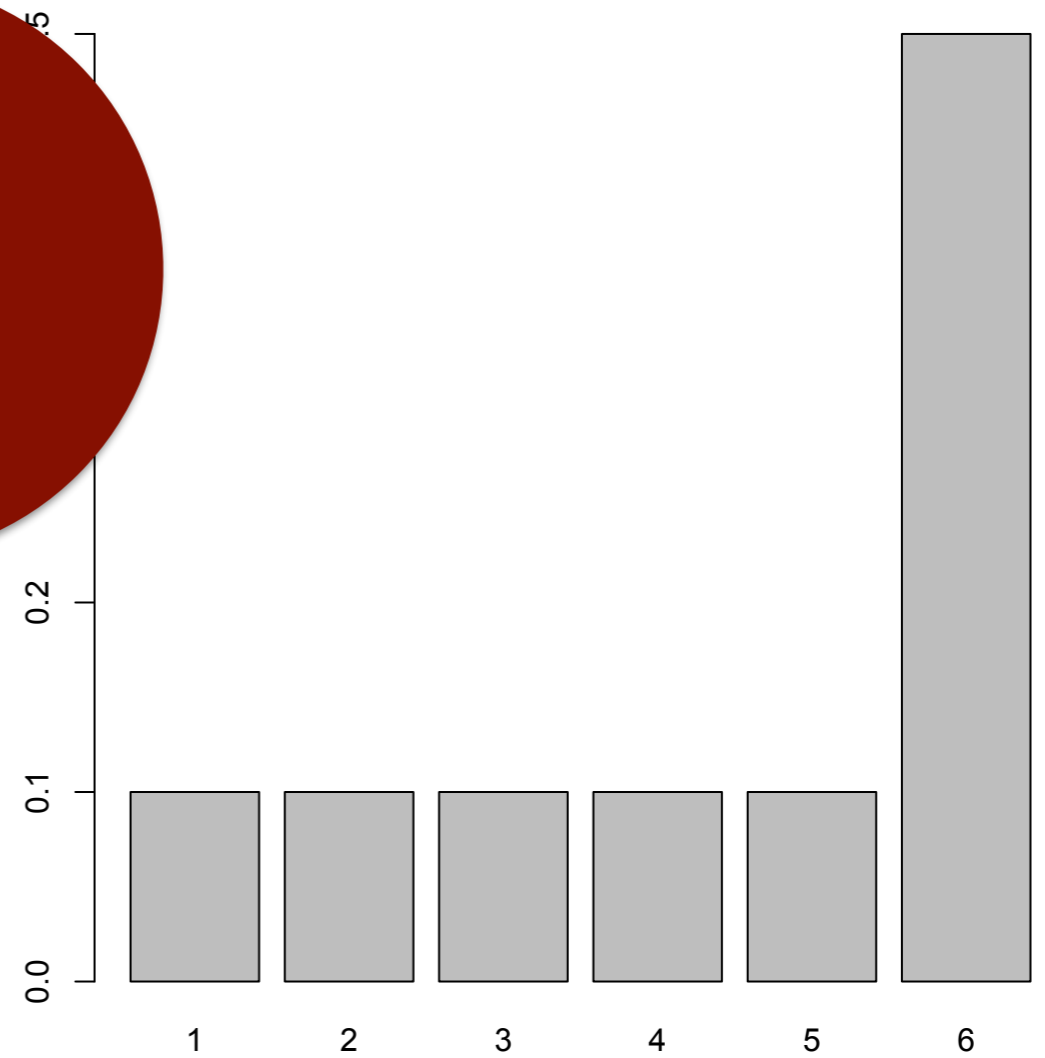
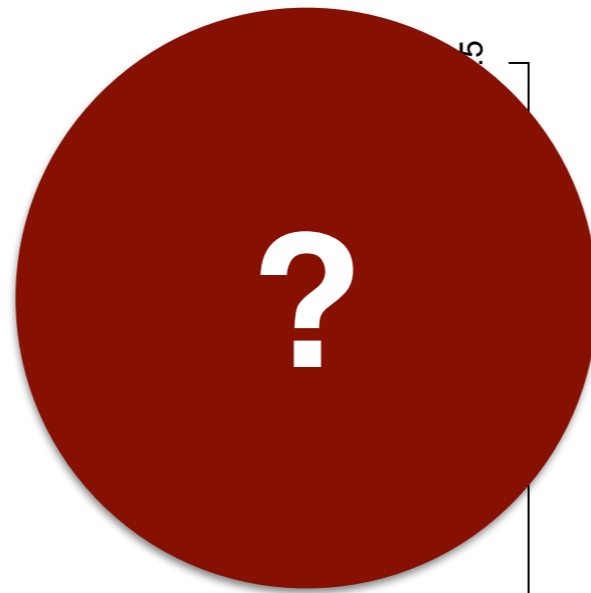
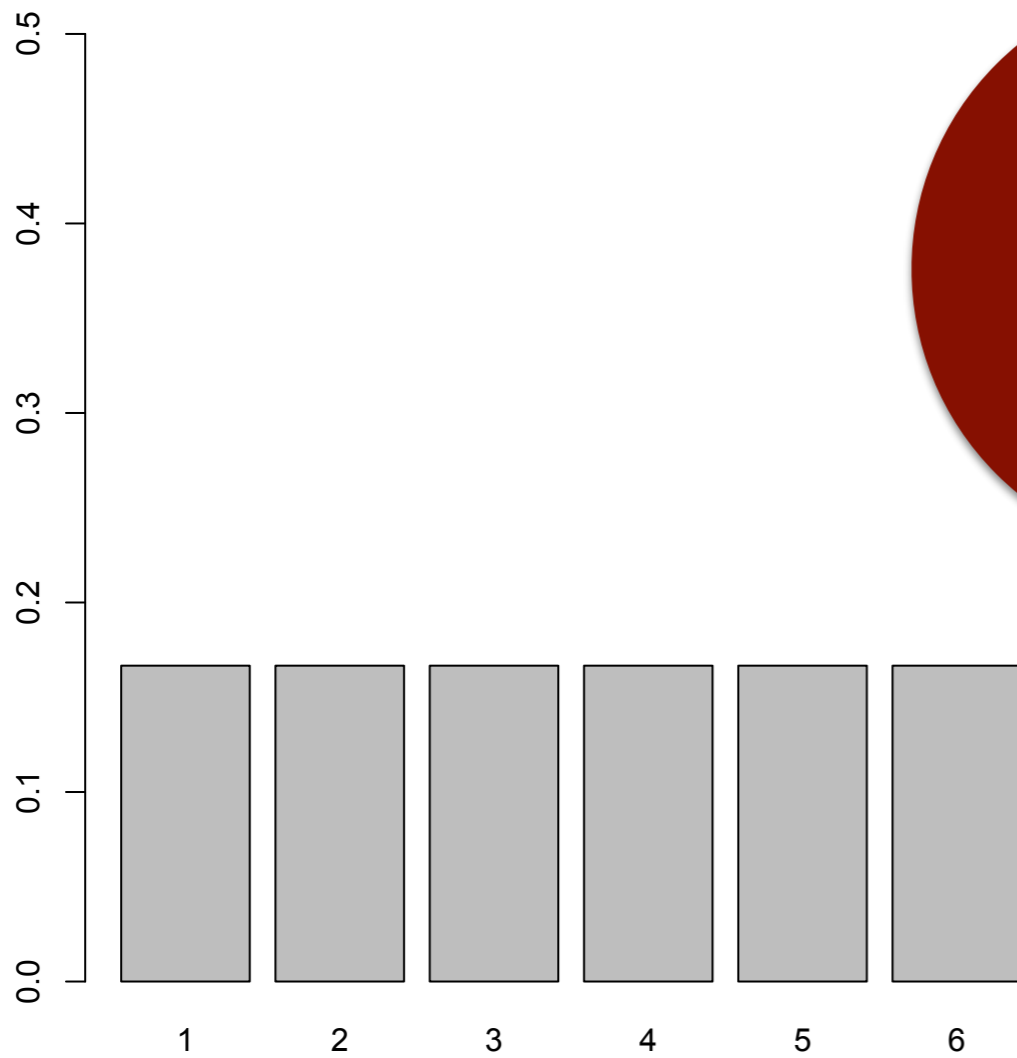


Probability



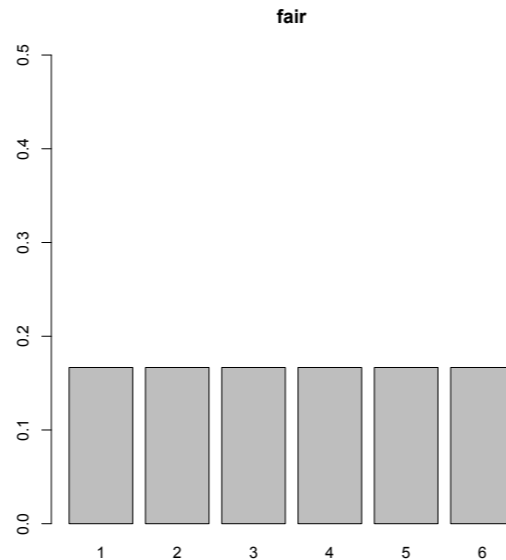
fair

not fair



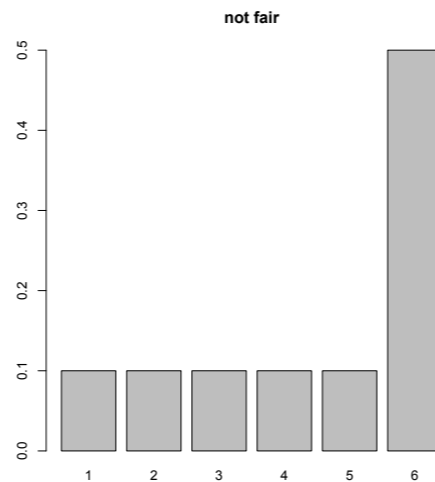
1. Data “Likelihood”

$$P(\text{2, 6, 6} \mid \text{fair})$$



$$= .17 \times .17 \times .17$$
$$= 0.004913$$

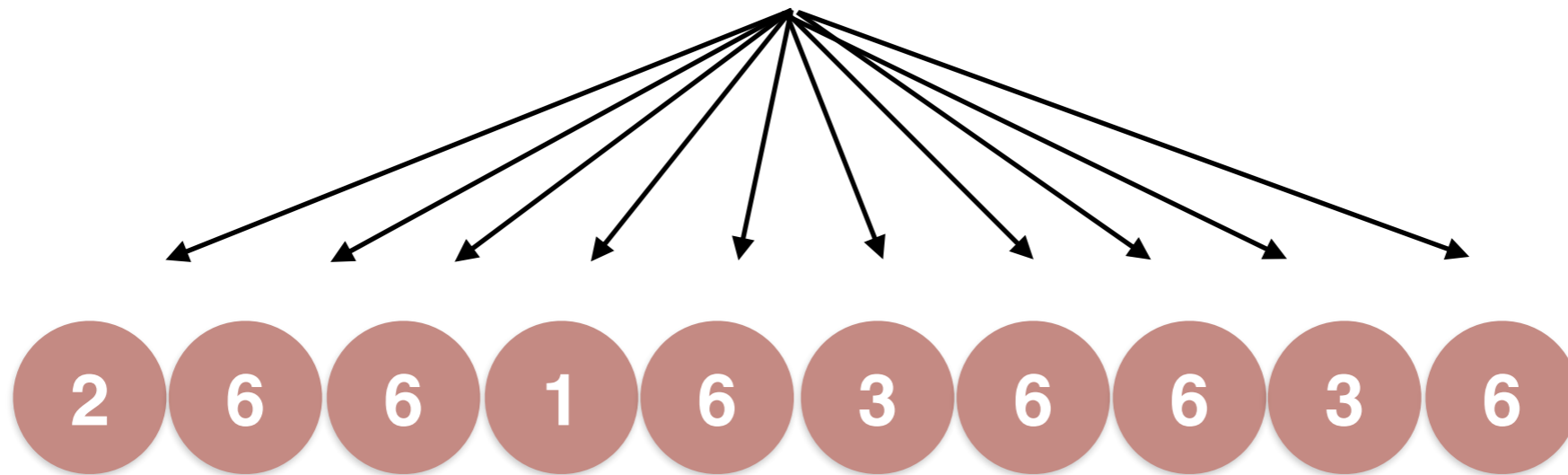
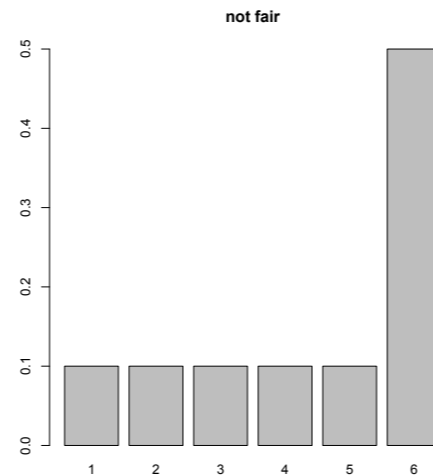
$$P(\text{2, 6, 6} \mid \text{not fair})$$



$$= .1 \times .5 \times .5$$
$$= 0.025$$

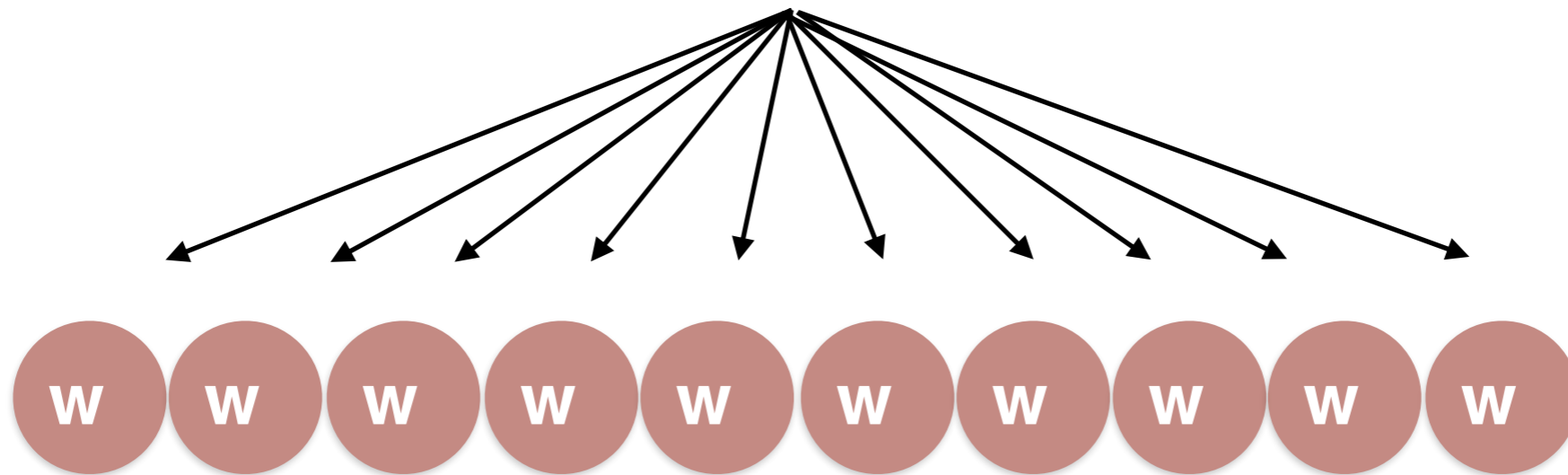
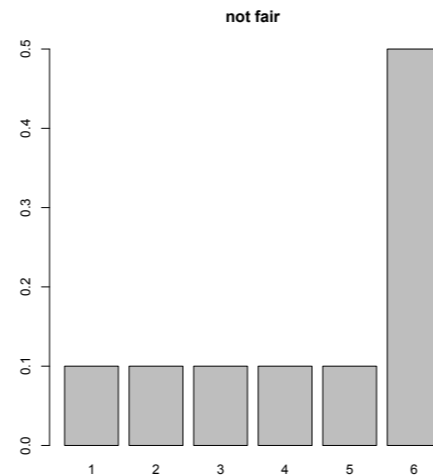
2. Conditional Probability

$$P(w \mid \theta)$$



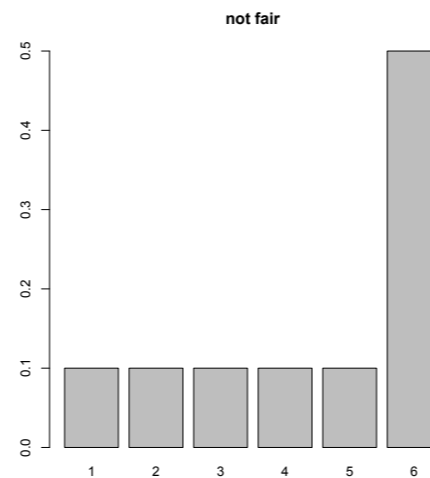
2. Conditional Probability

$$P(w \mid \theta)$$



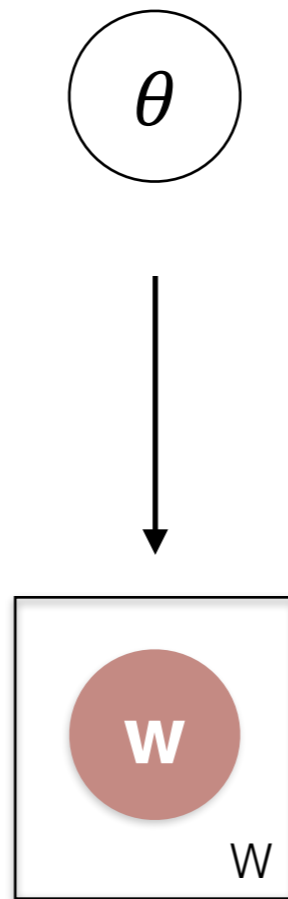
2. Conditional Probability

$$P(w \mid \theta)$$



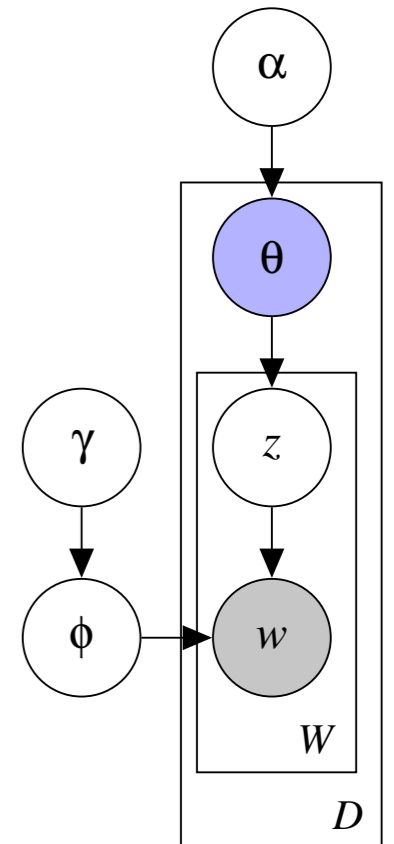
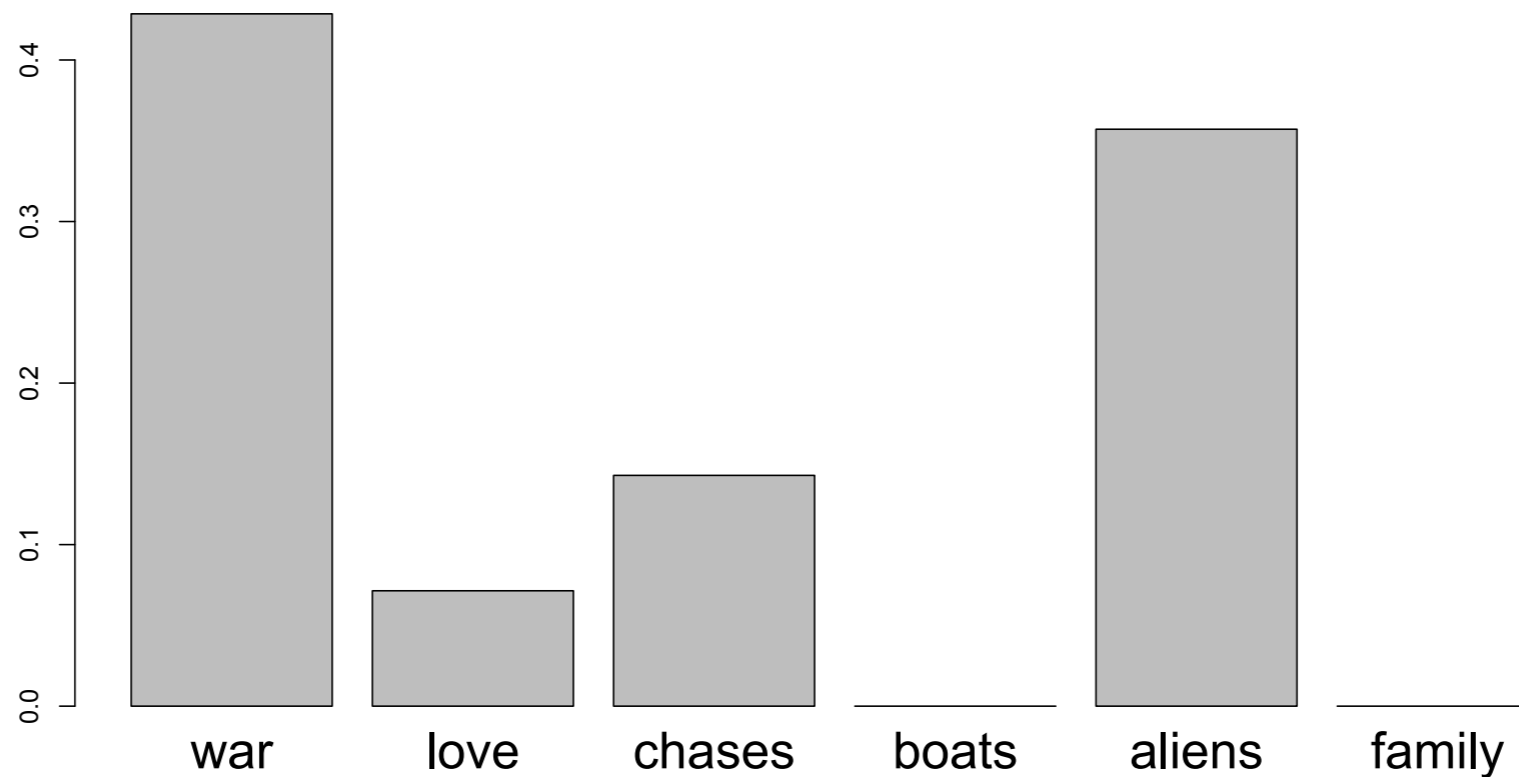
2. Conditional Probability

$$P(w \mid \theta)$$



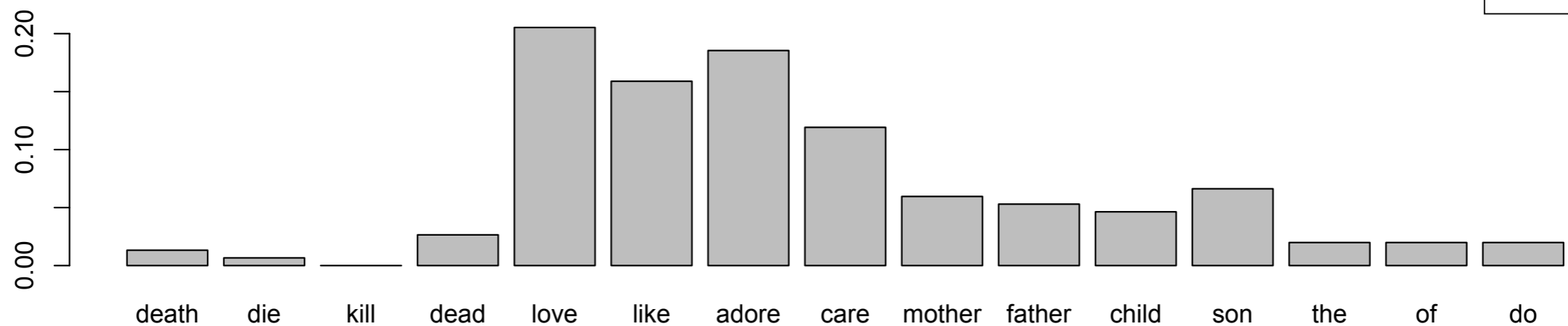
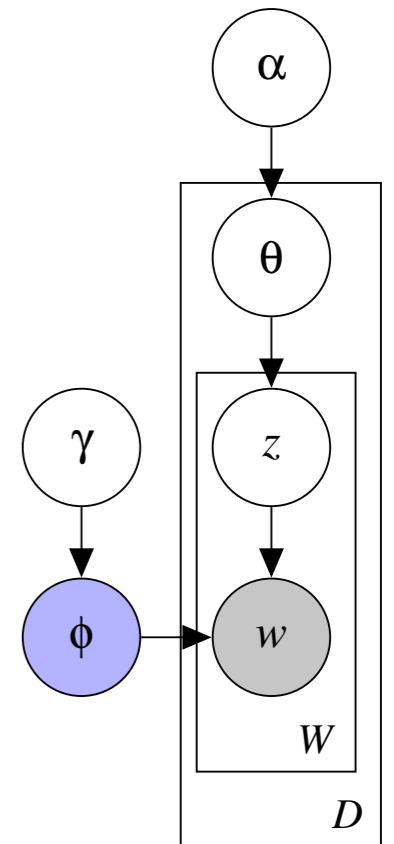
Topic Models

- A document has *distribution over topics*

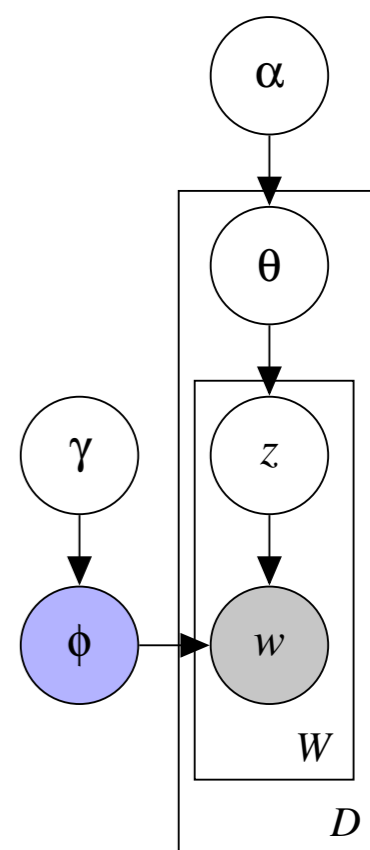
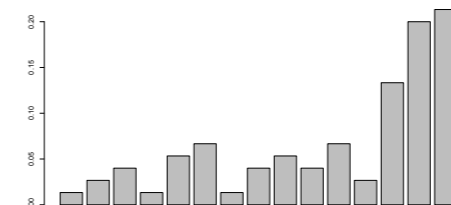
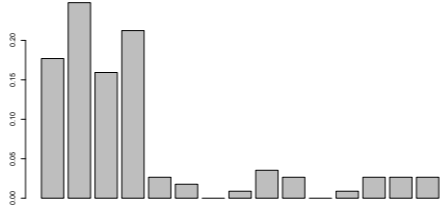
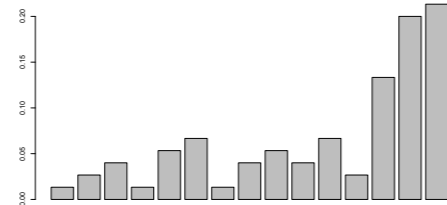
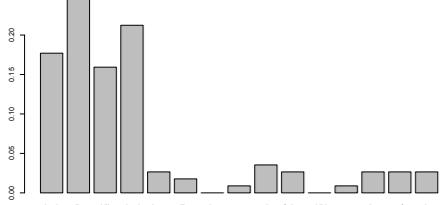
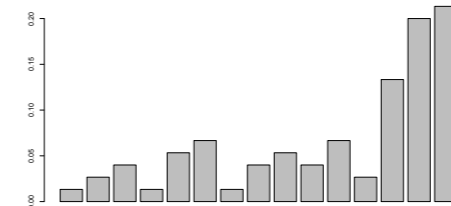
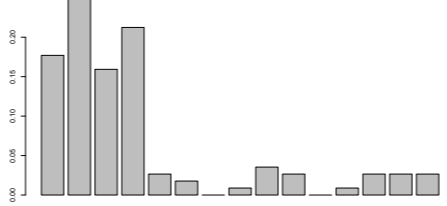
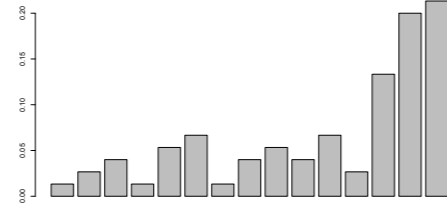
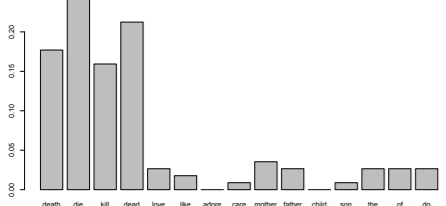
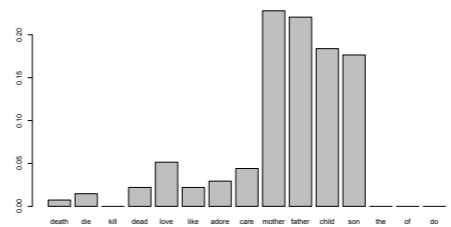
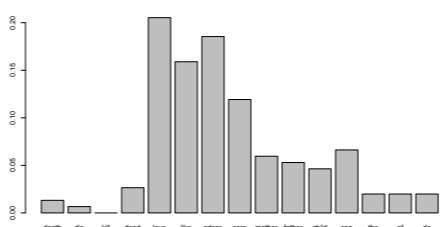
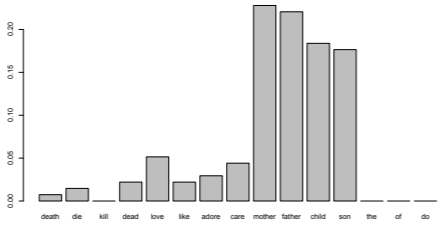
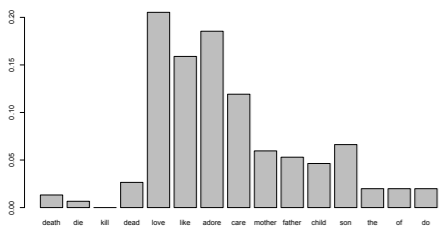
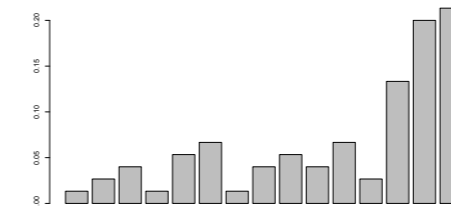
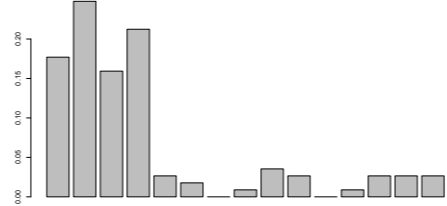
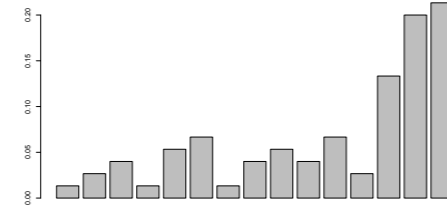
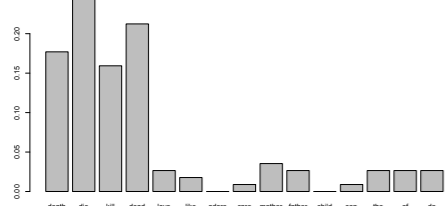
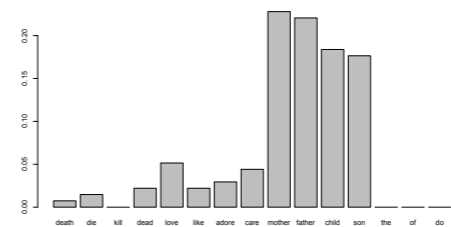
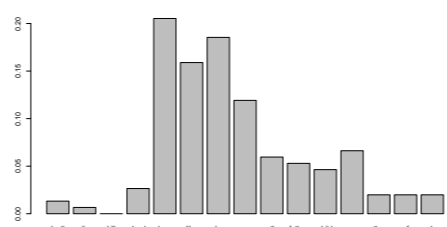
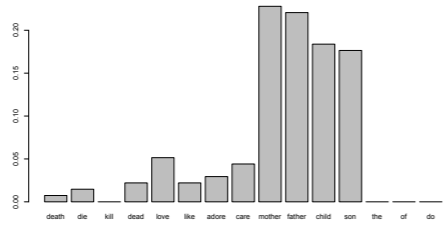
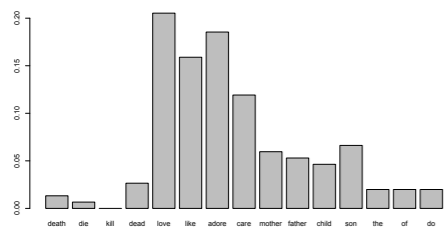


Topic Models

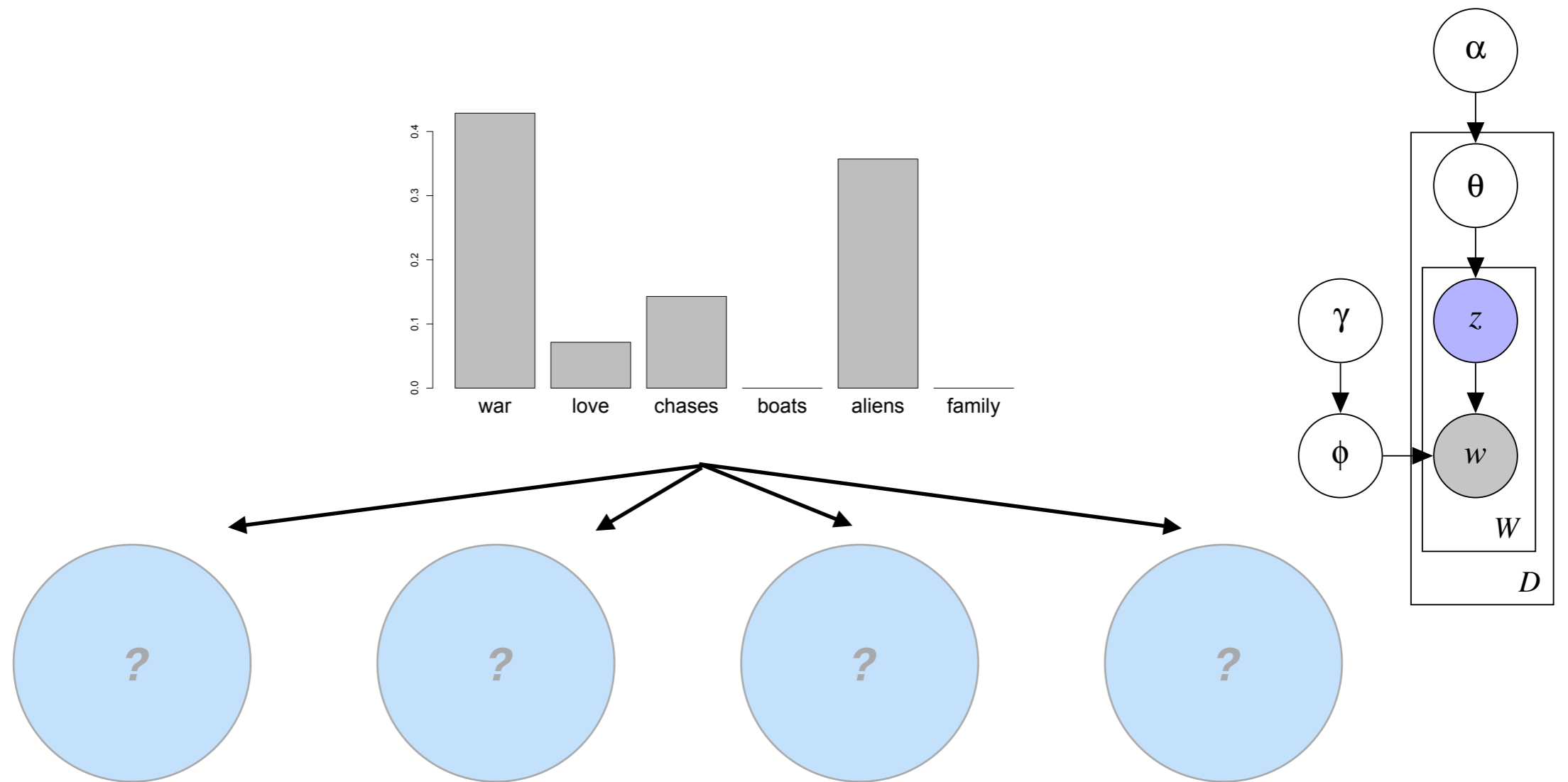
- A topic is a distribution over words



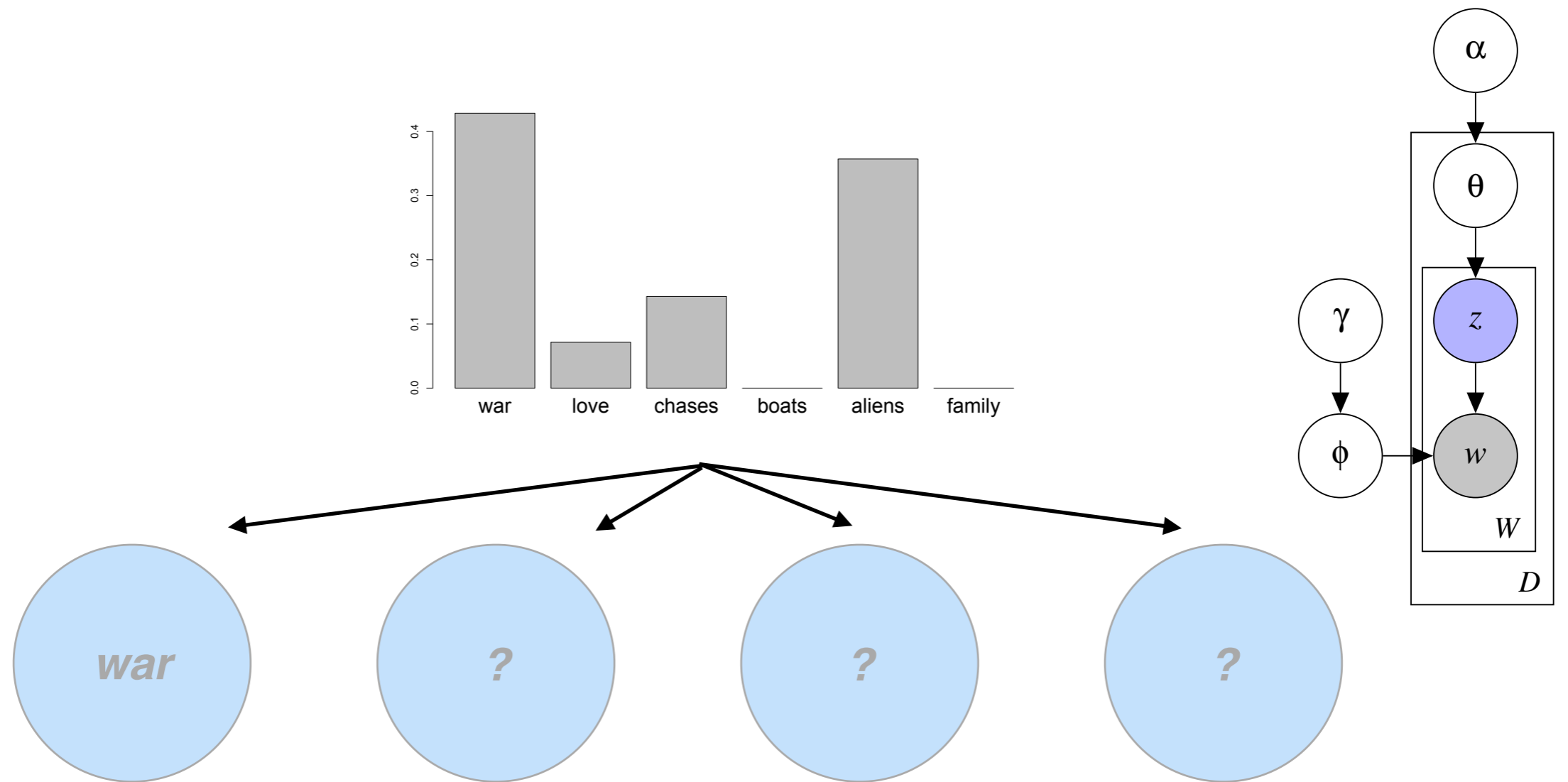
- e.g., $P(\text{"adore"} \mid \text{topic} = \text{love}) = .18$



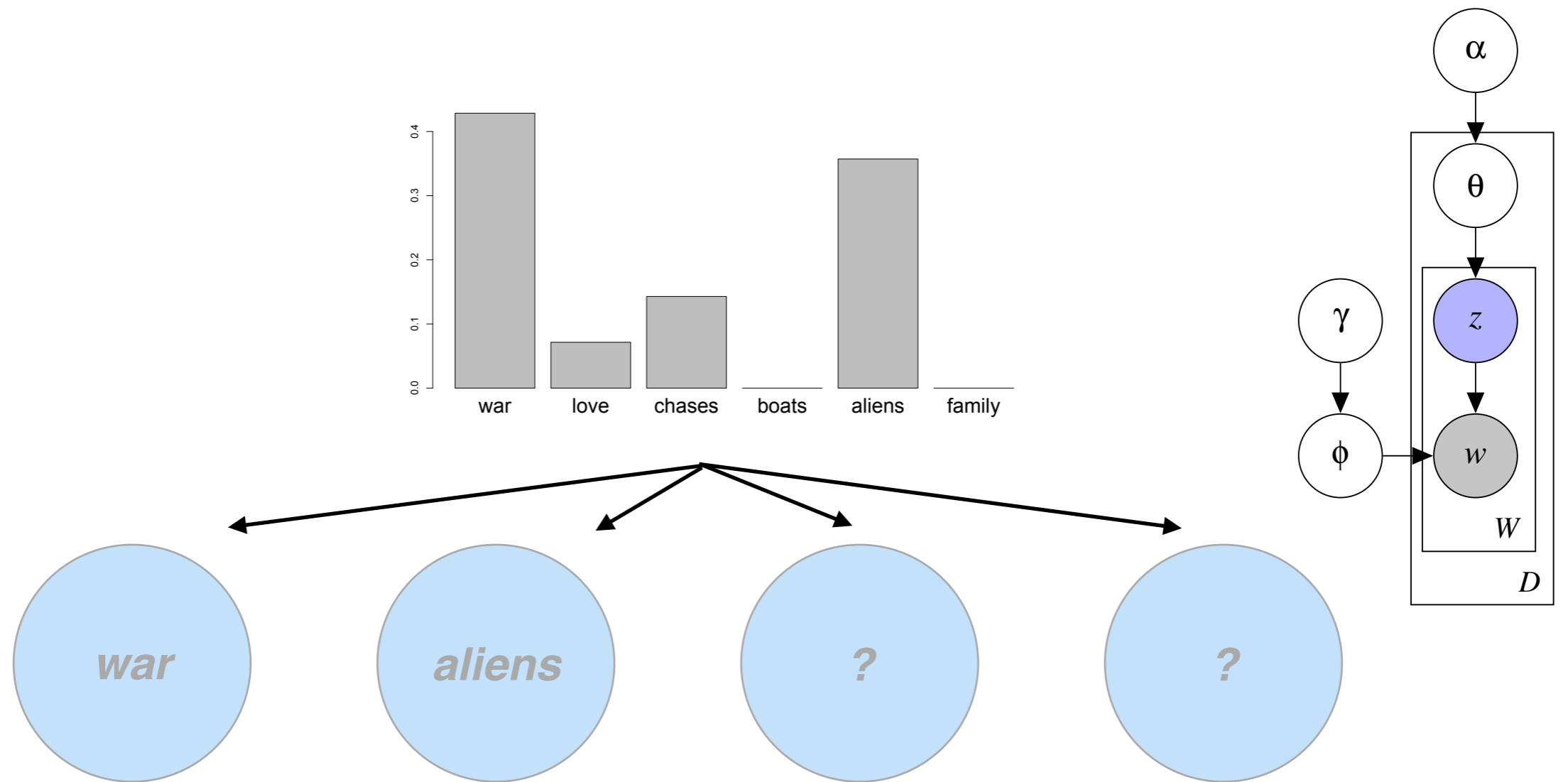
K=20



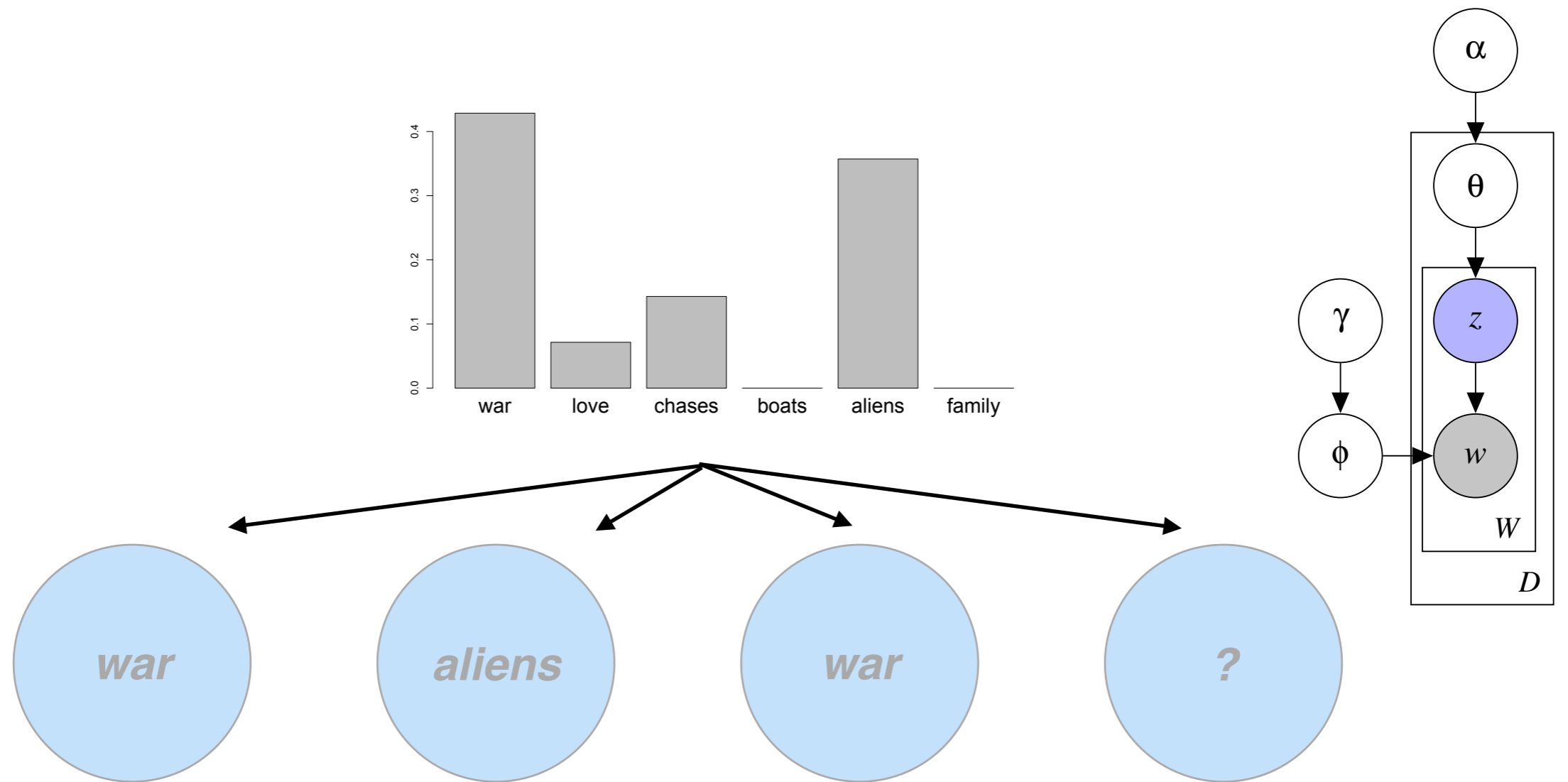
$P(\text{topic} \mid \text{topic distribution})$



$P(\text{topic} \mid \text{topic distribution})$



$P(\text{topic} \mid \text{topic distribution})$

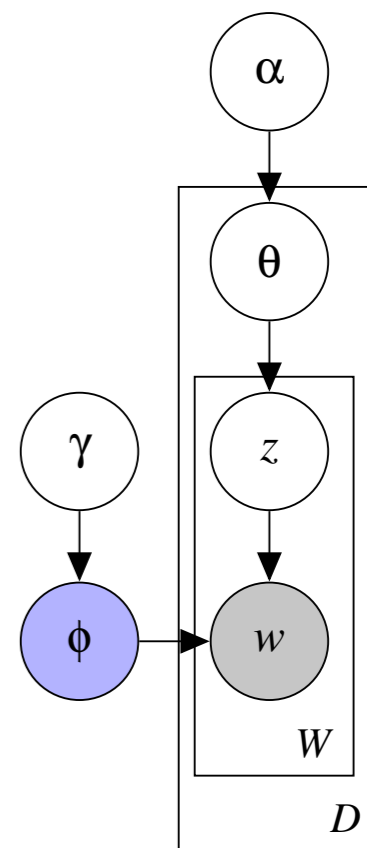
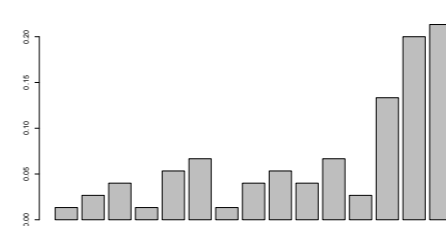
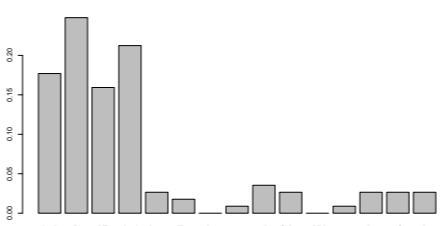
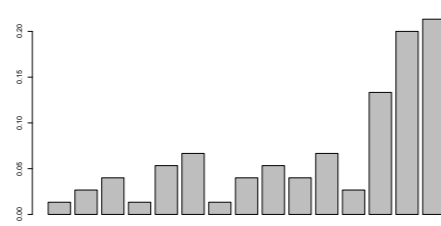
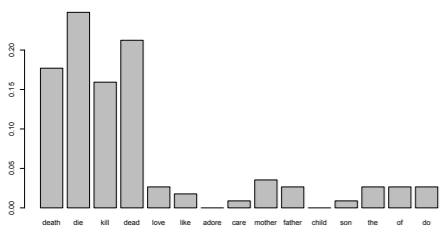
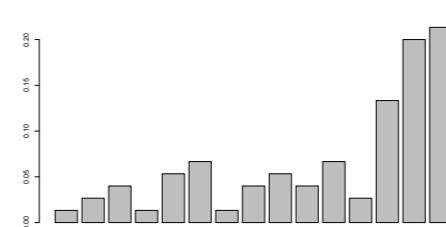
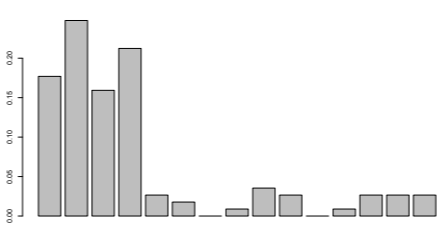
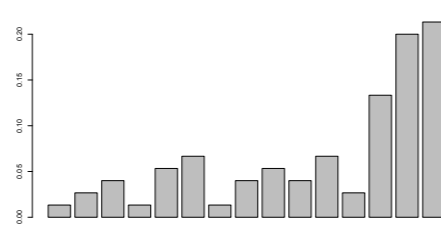
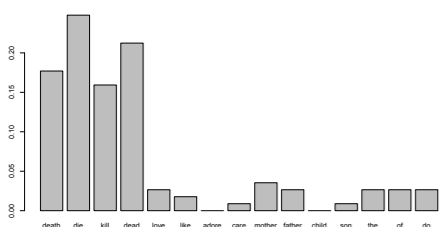
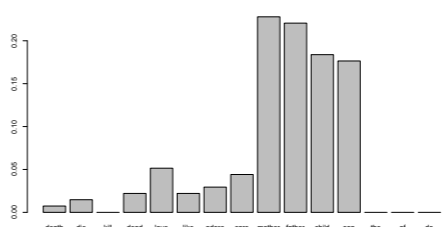
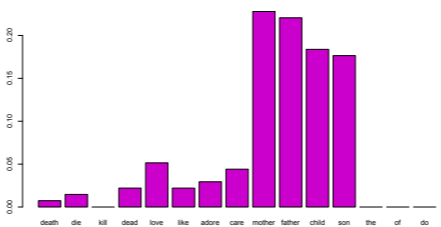
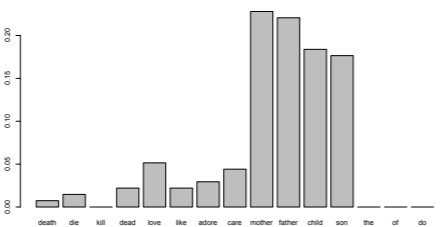
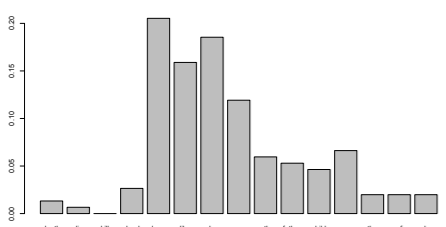
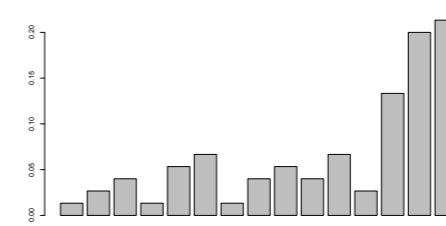
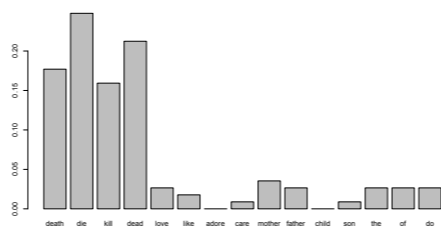
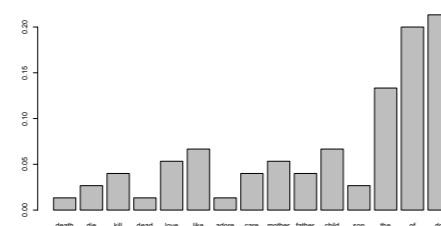
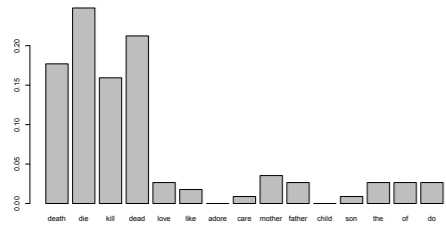
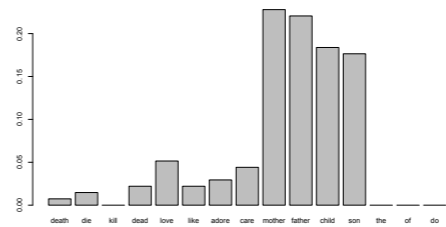
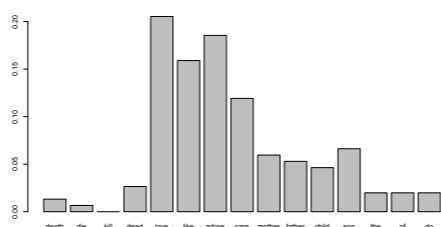
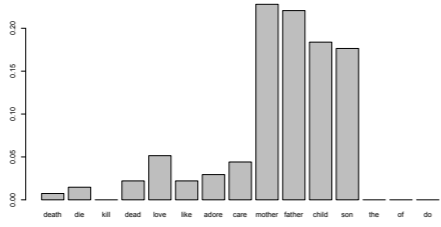
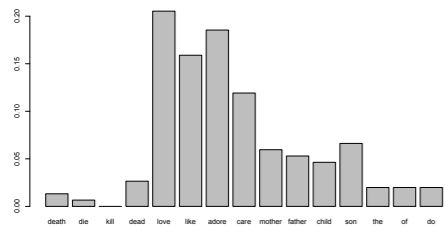


$P(\text{topic} \mid \text{topic distribution})$

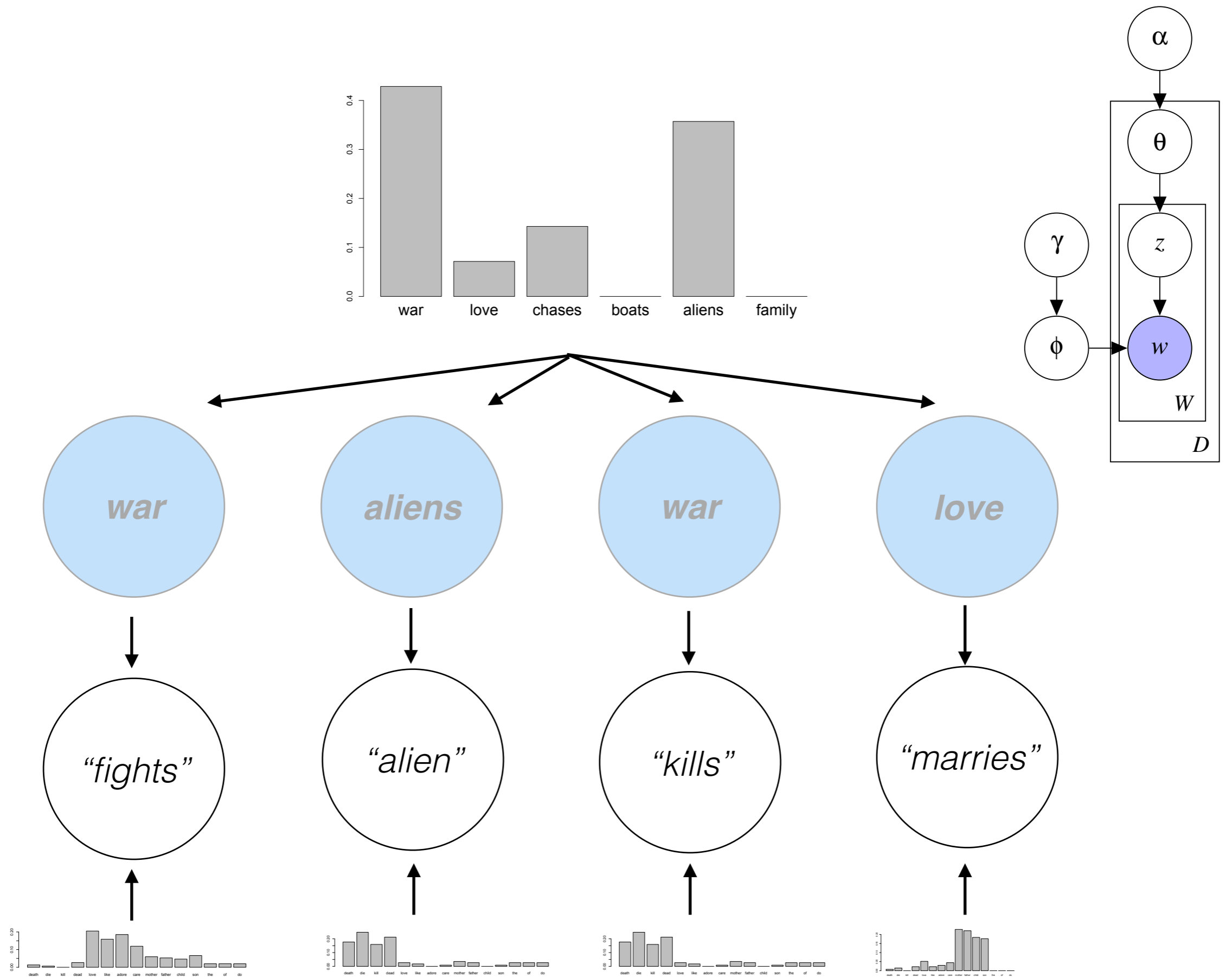


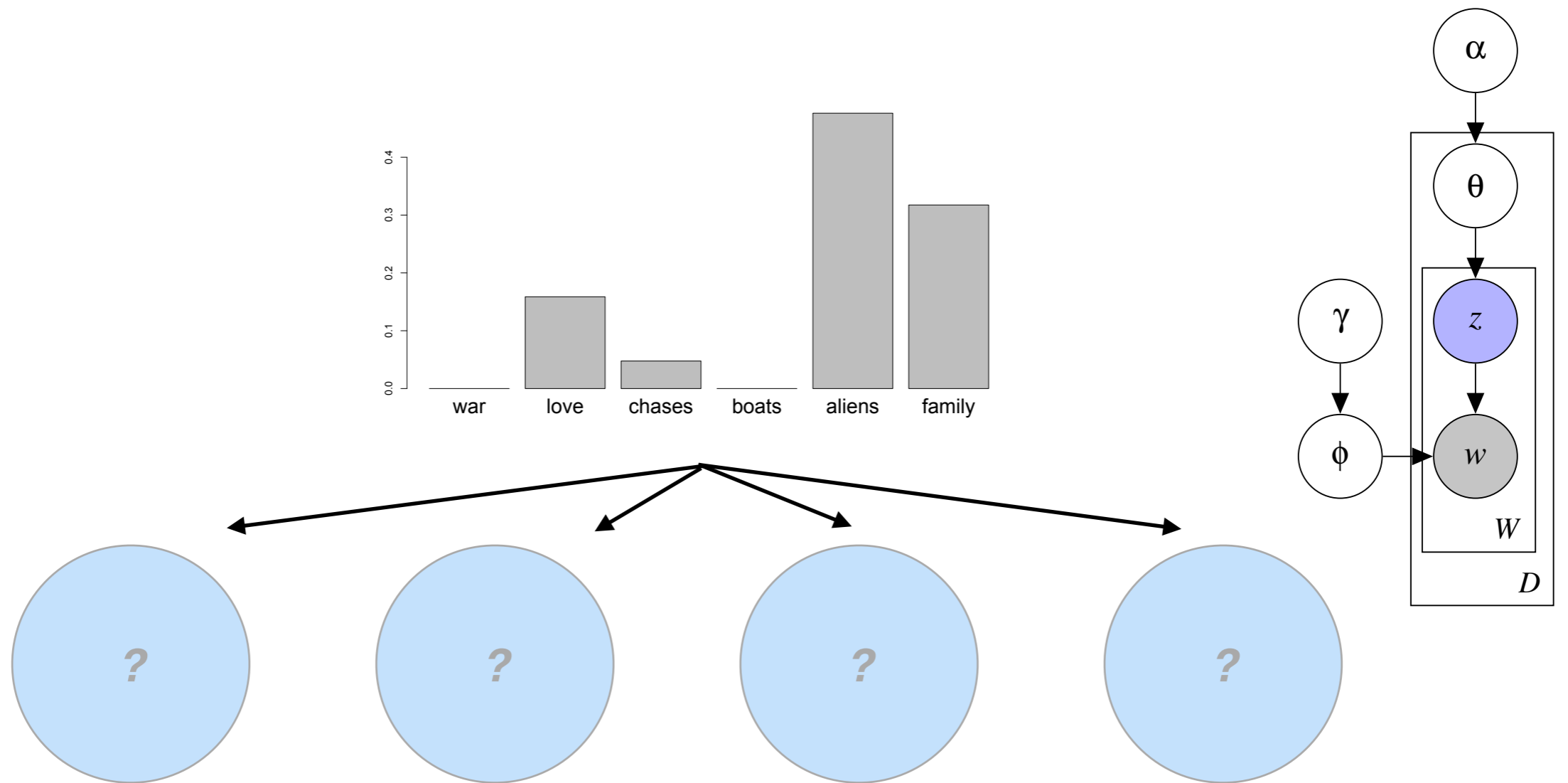
$P(\text{topic} \mid \text{topic distribution})$



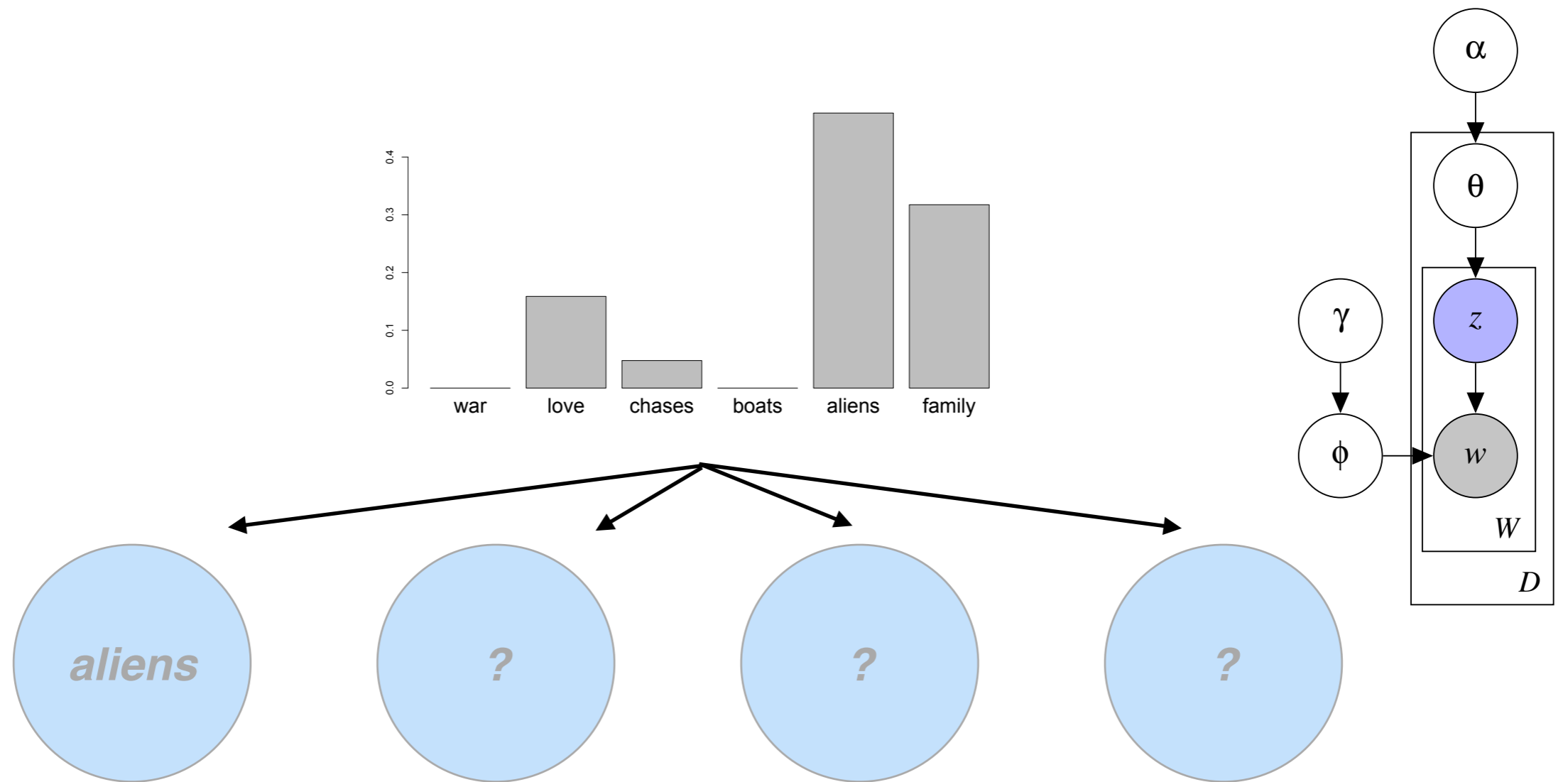


K=20

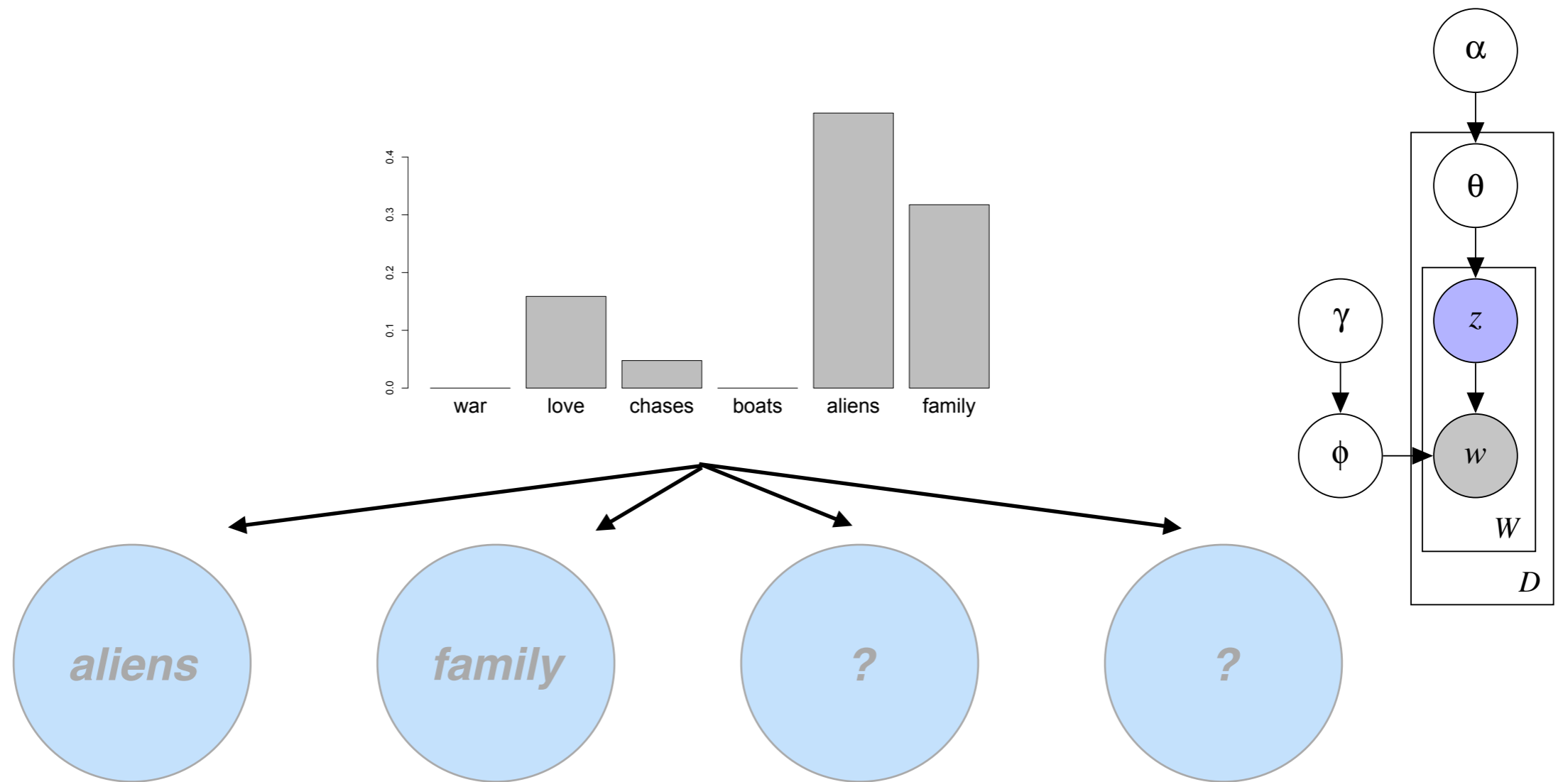




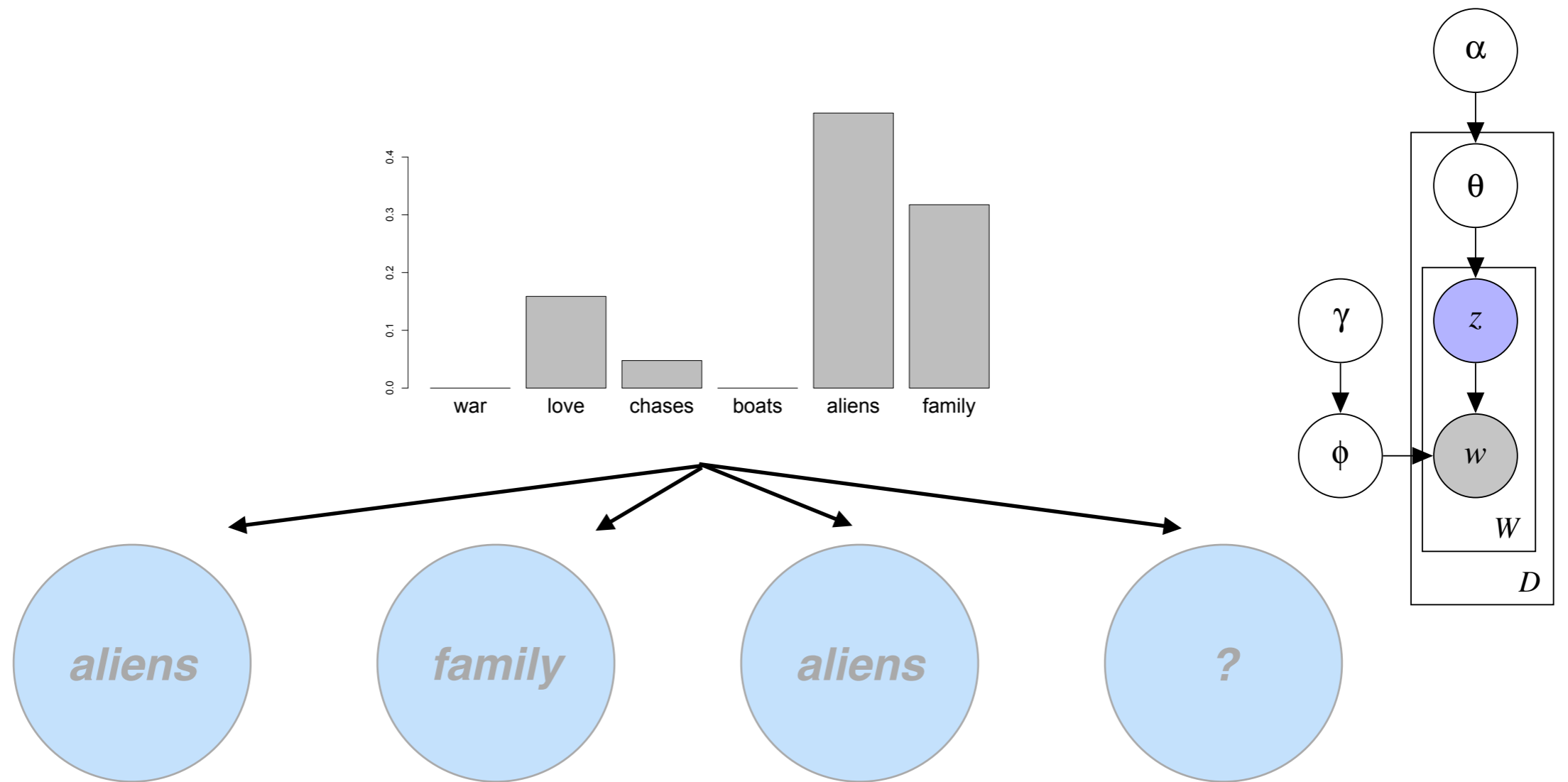
$P(\text{topic} \mid \text{topic distribution})$



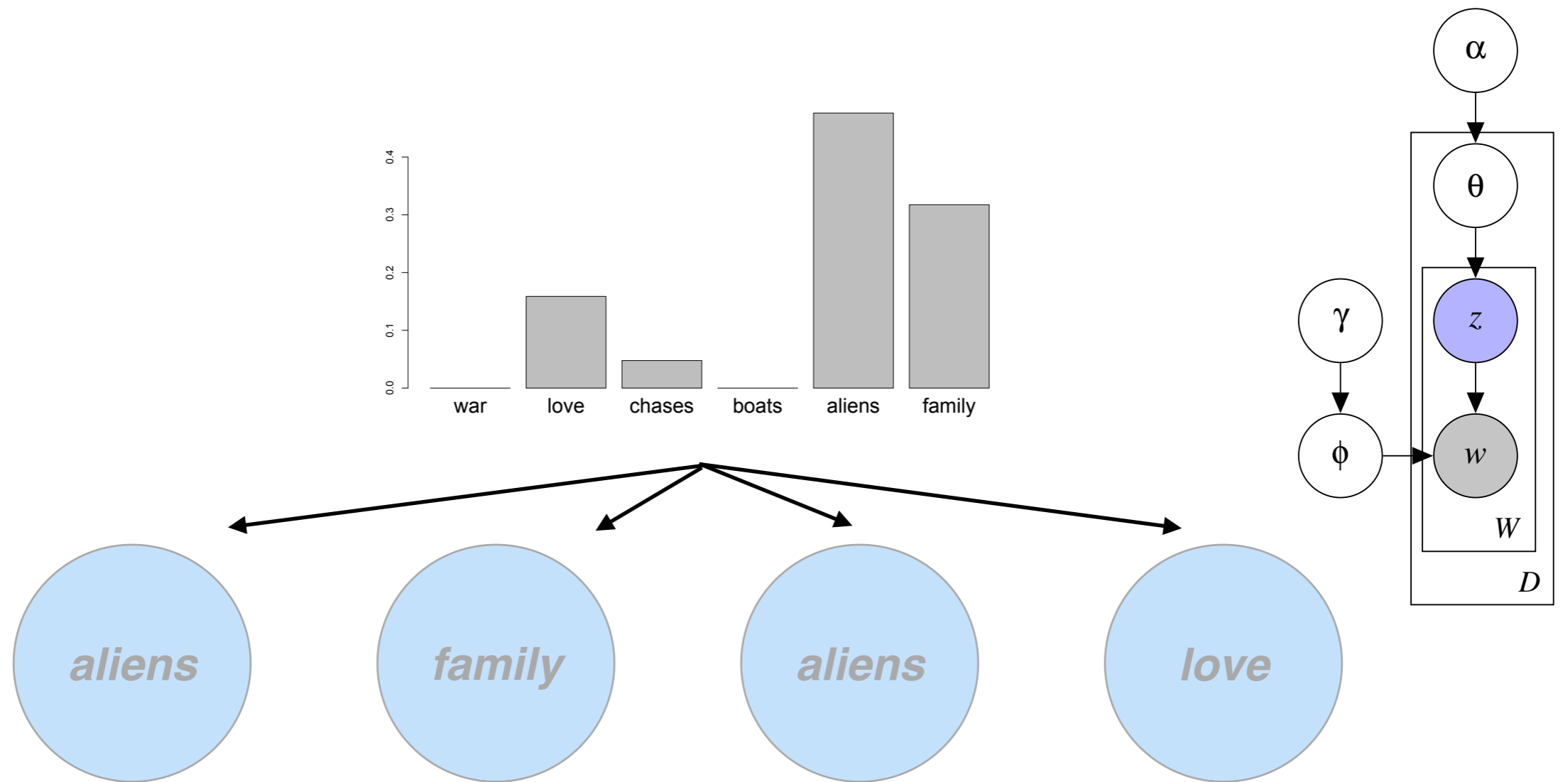
$P(\text{topic} \mid \text{topic distribution})$



$P(\text{topic} \mid \text{topic distribution})$

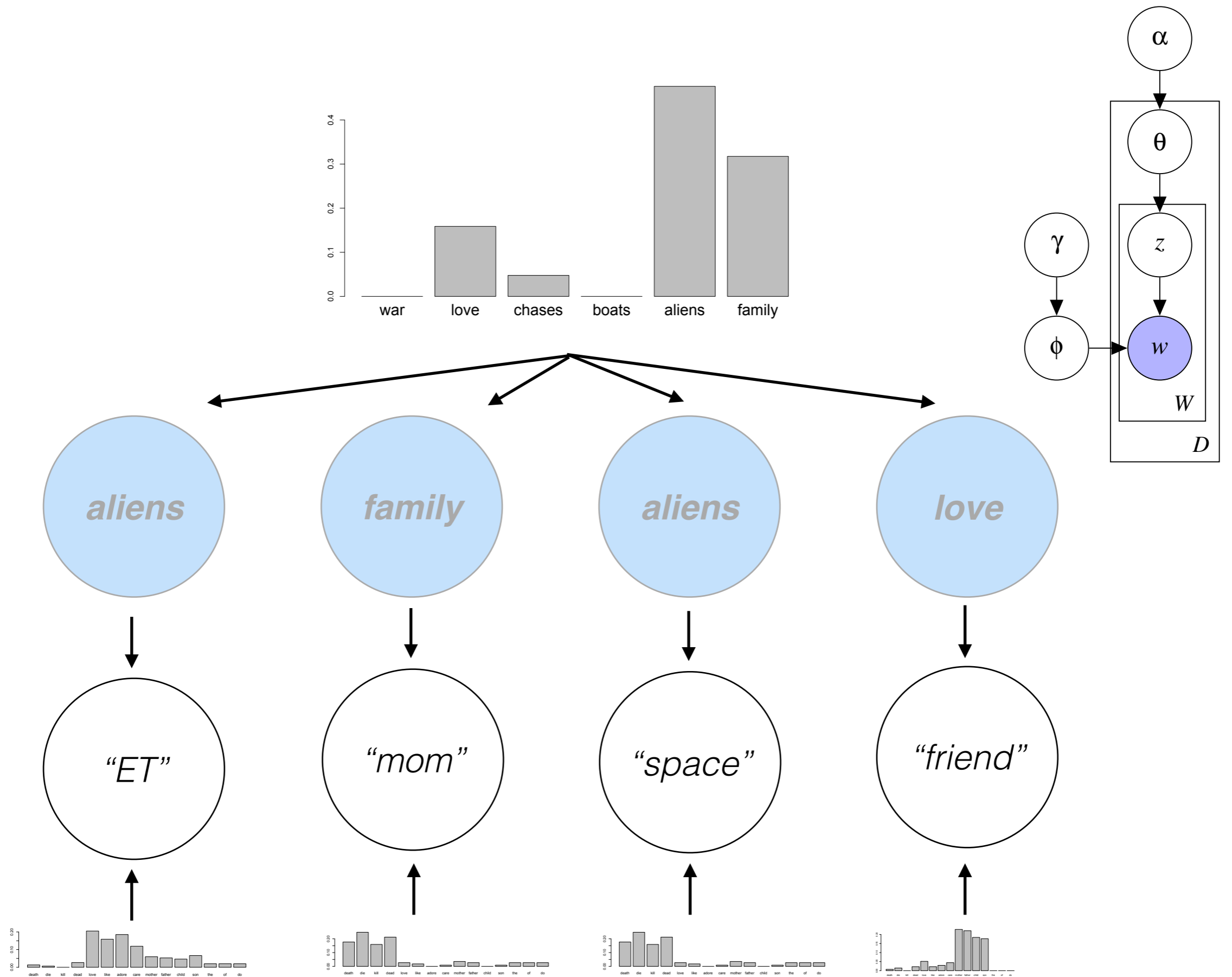


$P(\text{topic} \mid \text{topic distribution})$

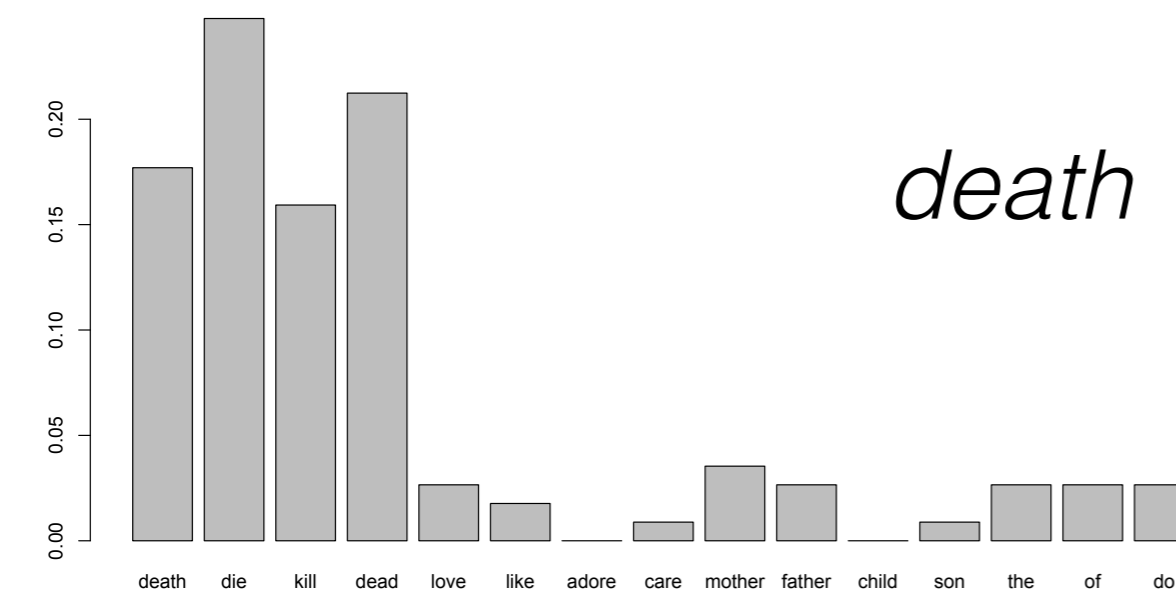
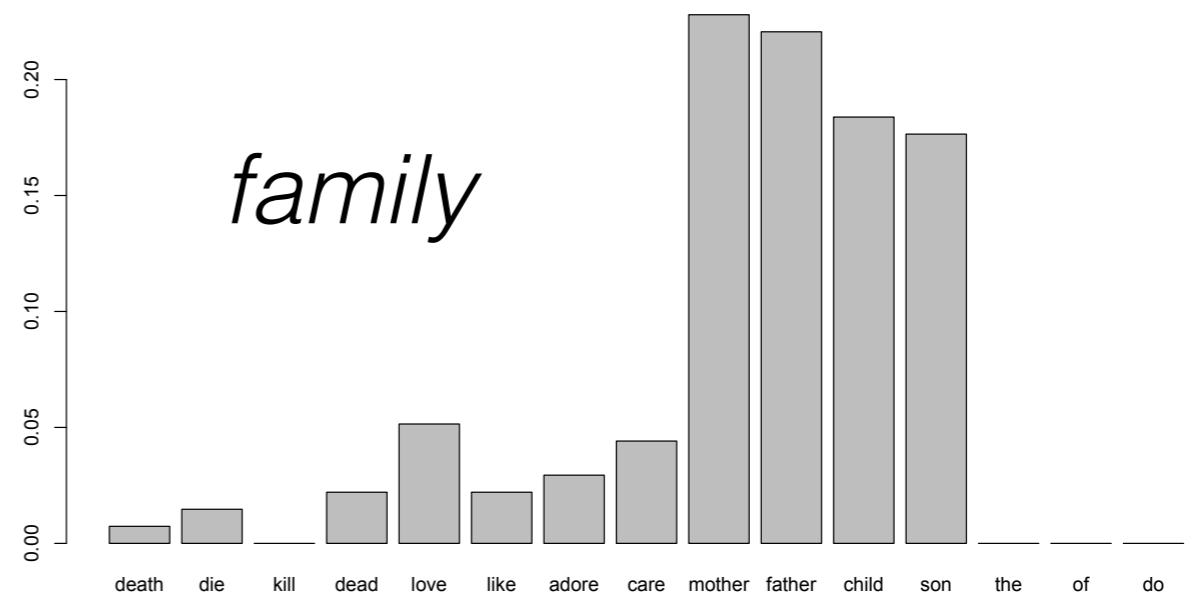
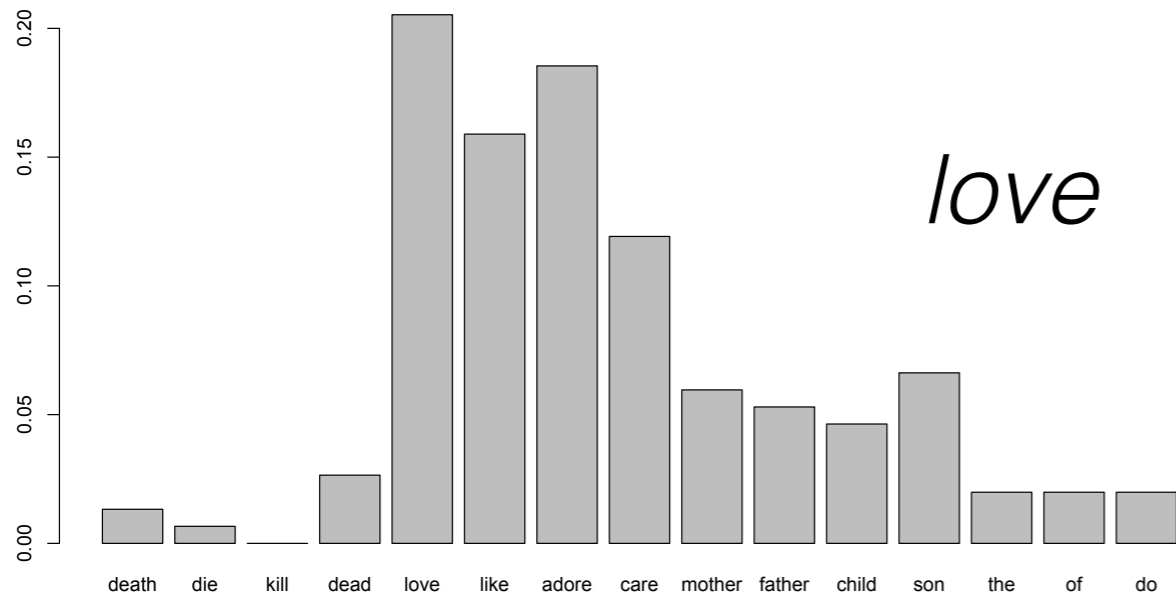


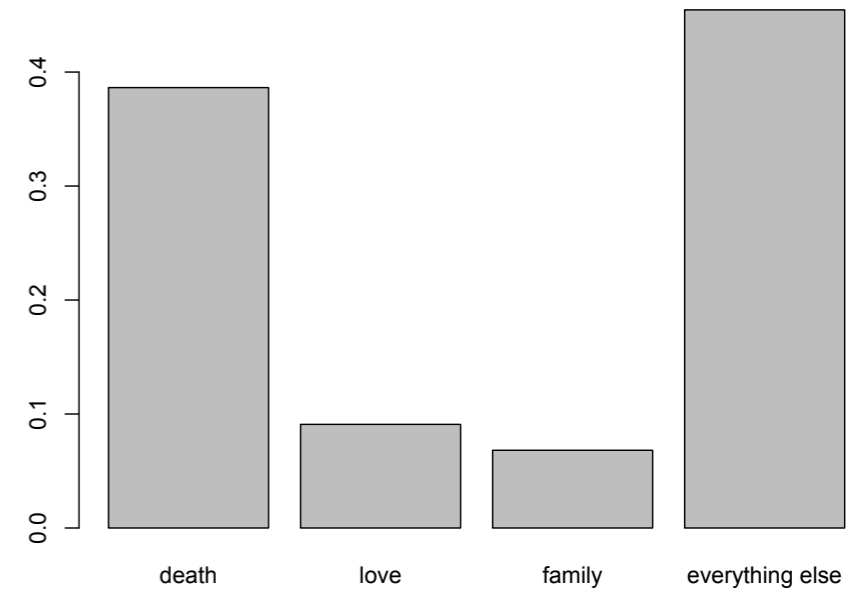
$P(\text{topic} \mid \text{topic distribution})$





Romeo and Juliet

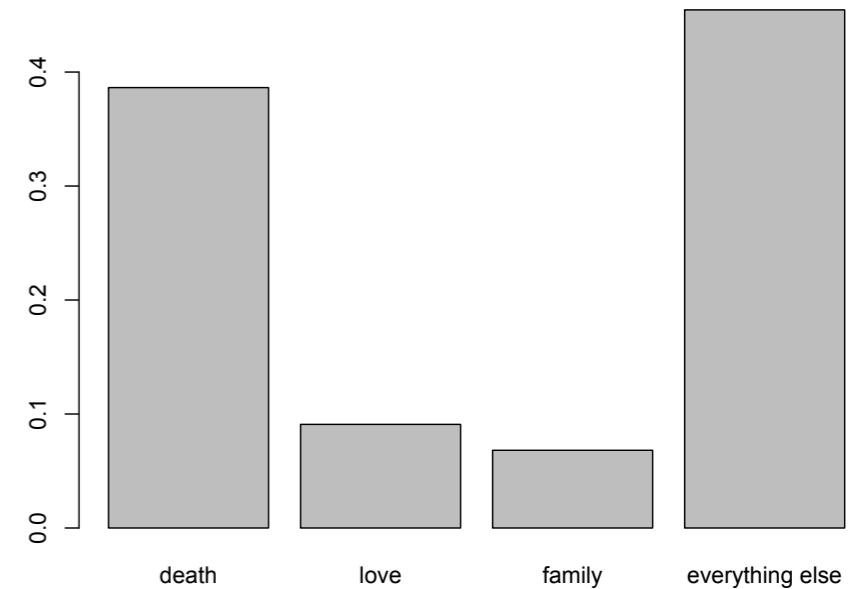




... The messenger, however, does not reach Romeo and, instead, Romeo learns of Juliet's apparent **death** from his servant Balthasar. Heartbroken, Romeo buys poison from an apothecary and goes to the Capulet **crypt**. He encounters Paris who has come to **mourn** Juliet privately. Believing Romeo to be a vandal, Paris **confronts** him and, in the ensuing **battle**, Romeo **kills** Paris. Still believing Juliet to be **dead**, he drinks the **poison**. Juliet then awakens and, finding Romeo **dead**, **stabs** herself with his **dagger**. The **feuding** families and the Prince meet at the **tomb** to find all three **dead**. Friar Laurence recounts the story of the two "star-cross'd lovers". The families are reconciled by their children's **deaths** and agree to end their **violent feud**. The play ends with the Prince's **elegy** for the lovers: "For never was a story of more woe / Than this of Juliet and her Romeo."

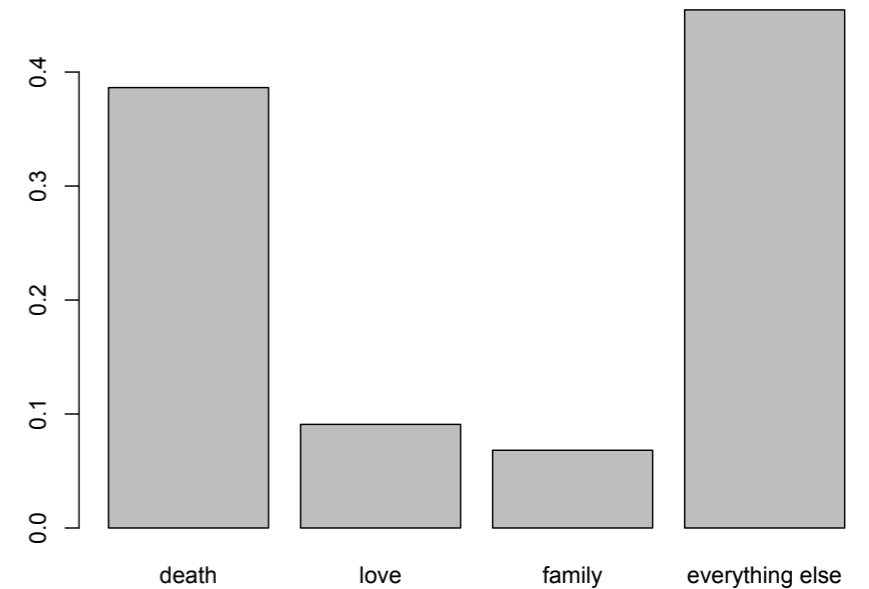
- **DEATH**
- LOVE
- FAMILY
- (EVERYTHING ELSE)

... The messenger, however, does not reach Romeo and, instead, Romeo learns of Juliet's apparent death from his servant Balthasar. Heartbroken, Romeo buys poison from an apothecary and goes to the Capulet crypt. He encounters Paris who has come to mourn Juliet privately. Believing Romeo to be a vandal, Paris confronts him and, in the ensuing battle, Romeo kills Paris. Still believing Juliet to be dead, he drinks the poison. Juliet then awakens and, finding Romeo dead, stabs herself with his dagger. The feuding families and the Prince meet at the tomb to find all three dead. Friar Laurence recounts the story of the two "star-cross'd lovers". The families are reconciled by their children's deaths and agree to end their violent feud. The play ends with the Prince's elegy for the lovers: "For never was a story of more woe / Than this of Juliet and her Romeo."

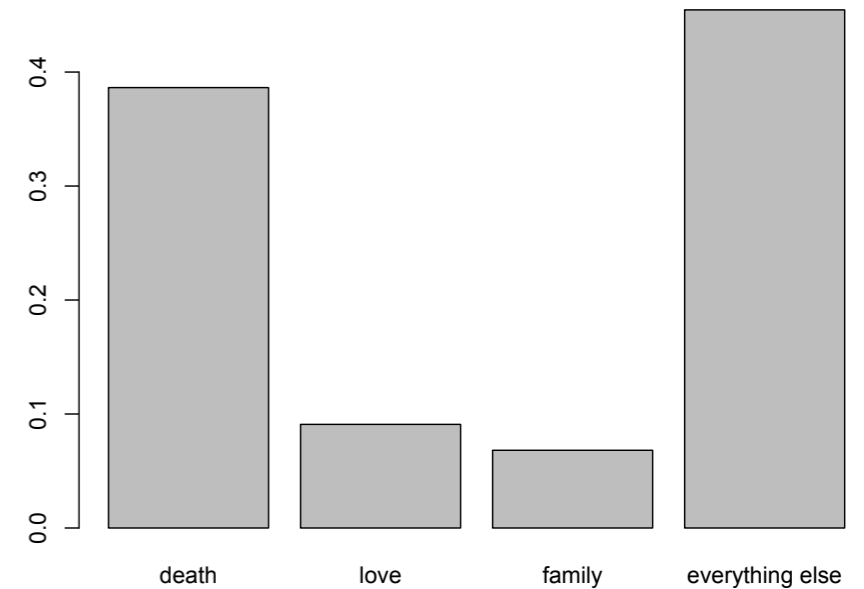


- DEATH
- LOVE
- FAMILY
- (EVERYTHING ELSE)

... The messenger, however, does not reach Romeo and, instead, Romeo learns of Juliet's apparent death from his servant Balthasar. Heartbroken, Romeo buys poison from an apothecary and goes to the Capulet crypt. He encounters Paris who has come to mourn Juliet privately. Believing Romeo to be a vandal, Paris confronts him and, in the ensuing battle, Romeo kills Paris. Still believing Juliet to be dead, he drinks the poison. Juliet then awakens and, finding Romeo dead, stabs herself with his dagger. The feuding **families** and the Prince meet at the tomb to find all three dead. Friar Laurence recounts the story of the two "star-cross'd lovers". The **families** are reconciled by their **children's** deaths and agree to end their violent feud. The play ends with the Prince's elegy for the lovers: "For never was a story of more woe / Than this of Juliet and her Romeo."



- DEATH
- LOVE
- **FAMILY**
- (EVERYTHING ELSE)



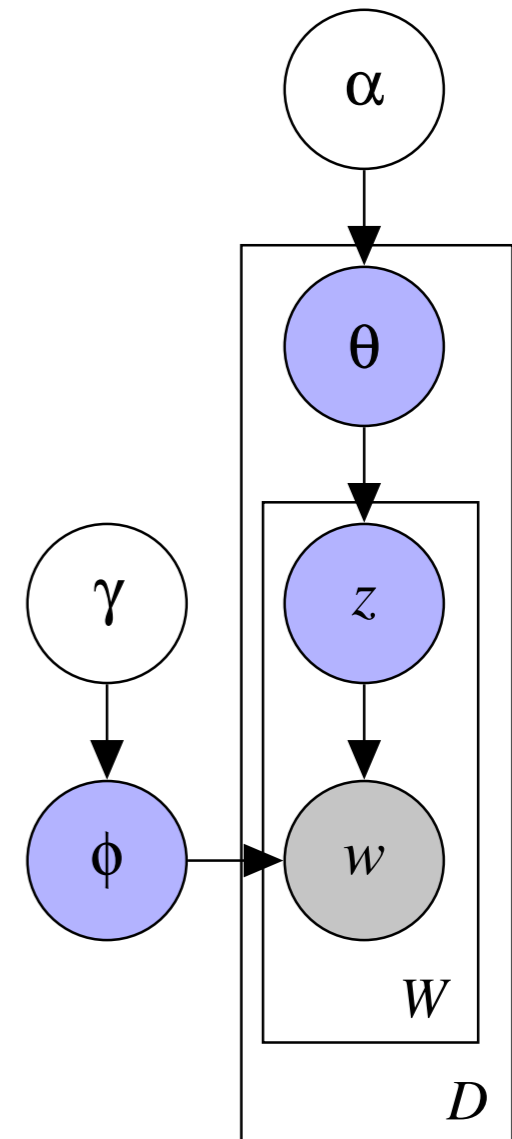
... The messenger, however, does not reach Romeo and, instead, Romeo learns of Juliet's apparent death from his servant Balthasar. Heartbroken, Romeo buys poison from an apothecary and goes to the Capulet crypt. He encounters Paris who has come to mourn Juliet privately. Believing Romeo to be a vandal, Paris confronts him and, in the ensuing battle, Romeo kills Paris. Still believing Juliet to be dead, he drinks the poison. Juliet then awakens and, finding Romeo dead, stabs herself with his dagger. The feuding families and the Prince meet at the tomb to find all three dead. Friar Laurence recounts the story of the two "star-cross'd lovers". The families are reconciled by their children's deaths and agree to end their violent feud. The play ends with the Prince's elegy for the lovers: "For never was a story of more woe / Than this of Juliet and her Romeo."

- DEATH
- LOVE
- FAMILY
- (EVERYTHING ELSE)

Inference

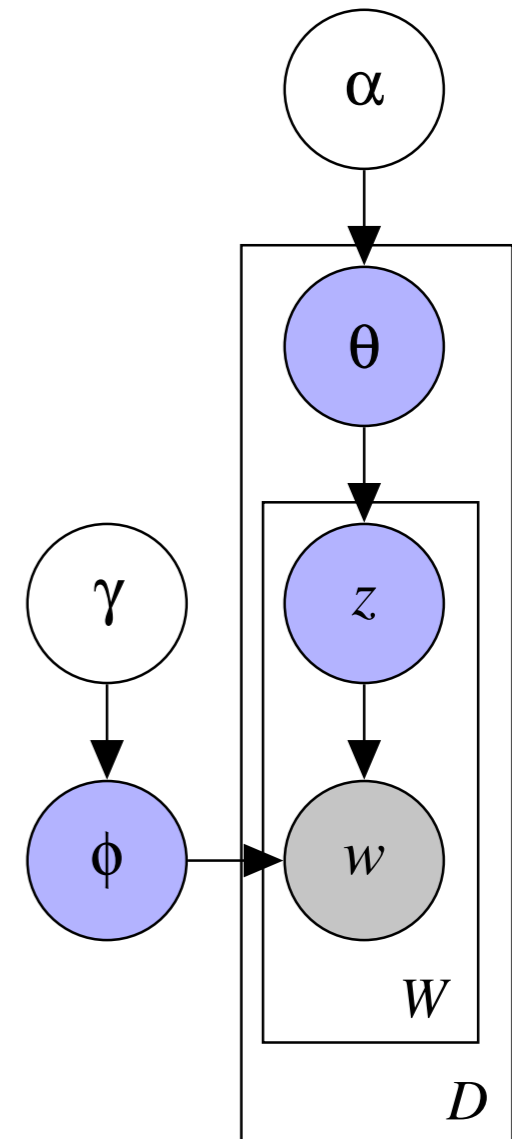
- What are the topic distributions for each document?
- What are the topic assignments for each word in a document?
- What are the word distributions for each topic?

Find the parameters that maximize the likelihood of the data!



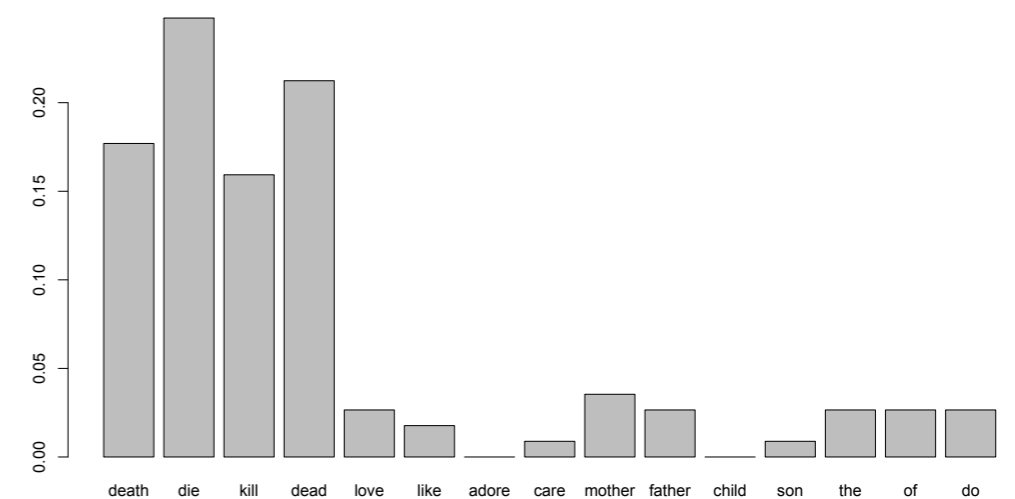
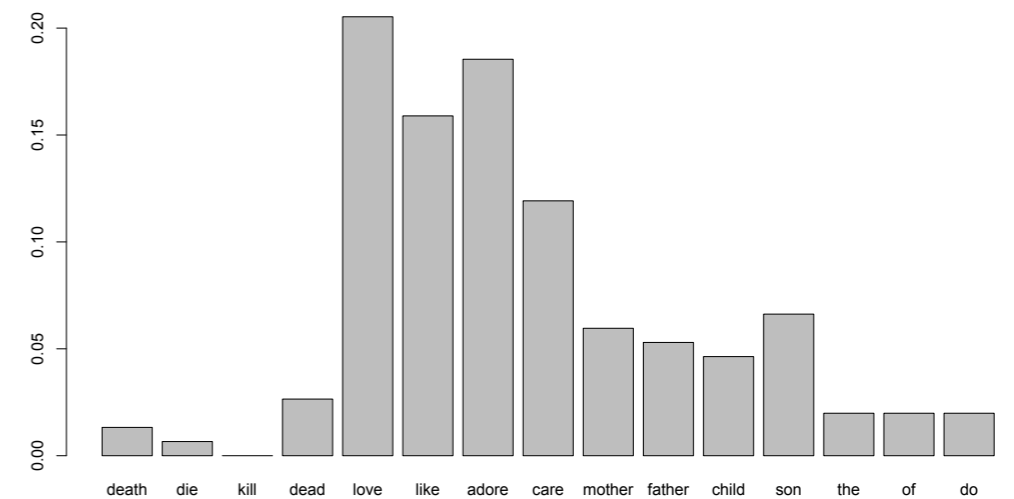
Gibbs Sampling

- 1. Start with some initial value for all the variables
- 2. Sample a value for a variable conditioned on all of the other variables around it (using Bayes' theorem)



Inferred Topics

{album, band, music}	{government, party, election}	{game, team, player}
album	government	game
band	party	team
music	election	player
song	state	win
release	political	play
{god, call, give}	{company, market, business}	{math, number, function}
god	company	math
call	market	number
give	business	function
man	year	code
time	product	set
{city, large, area}	{math, energy, light}	{law, state, case}
city	math	law
large	energy	state
area	light	case
station	field	court
include	star	legal



Examples

- Mining the Dispatch
<http://dsl.richmond.edu/dispatch/>
- Wikipedia Topics
<http://www.princeton.edu/~achaney/tmve/wiki100k/browse/topic-list.html>
- Quiet Transformations of Literary Studies
<http://www.rci.rutgers.edu/~ag978/quiet/>

Try it yourself

- book summaries, movie summaries, PMLA ,
Classical Quarterly, Renaissance Quarterly,
Shakespeare + English stoplist
<http://bit.ly/1hdKX0R>
- Topic Modeling Tool
<https://code.google.com/p/topic-modeling-tool/>

Representation Learning

Assume we've trained a **logistic regression** classifier to predict whether a tweet was written by a person who lives in Chicago.

β_{Chicago}

i	feat	value
1	I	0.004
2	live	0.0013
3	in	-0.001
4	New York	-13.7
5	Chicago	8.7
6	Boston	-10.8
7	Pittsburgh	-5.7
8	snow	2.7

Representation Learning

“I live in Chicago”

$\beta_{\text{Chicago}} =$

i	feat	value
1	I	0.004
2	live	0.0013
3	in	-0.01
4	New York	-13.7
5	Chicago	8.7
6	Boston	-10.8
7	Pittsburgh	-5.7
8	snow	2.7

$X =$

i	feat	value
1	I	1
2	live	1
3	in	1
4	New York	0
5	Chicago	1
6	Boston	0
7	Pittsburgh	0
8	snow	0

Representation Learning

“I live in Chicagoland”

$\beta_{\text{Chicago}} =$

i	feat	value
1	I	0.004
2	live	0.0013
3	in	-0.01
4	New York	-13.7
5	Chicago	8.7
6	Boston	-10.8
7	Pittsburgh	-5.7
8	snow	2.7

$X =$

i	feat	value
1	I	1
2	live	1
3	in	1
4	New York	0
5	Chicago	0
6	Boston	0
7	Pittsburgh	0
8	snow	0

Representation Learning

- Learn **alternate** representations for inputs (and sometimes outputs) aside from their raw (atomic) values.
- For words, this generally means representations that encode some measure of similarity.
 - Hard word clusters (e.g., Brown clusters)
 - Low-dimensional “embeddings” ($w \in \mathbb{R}^K$)

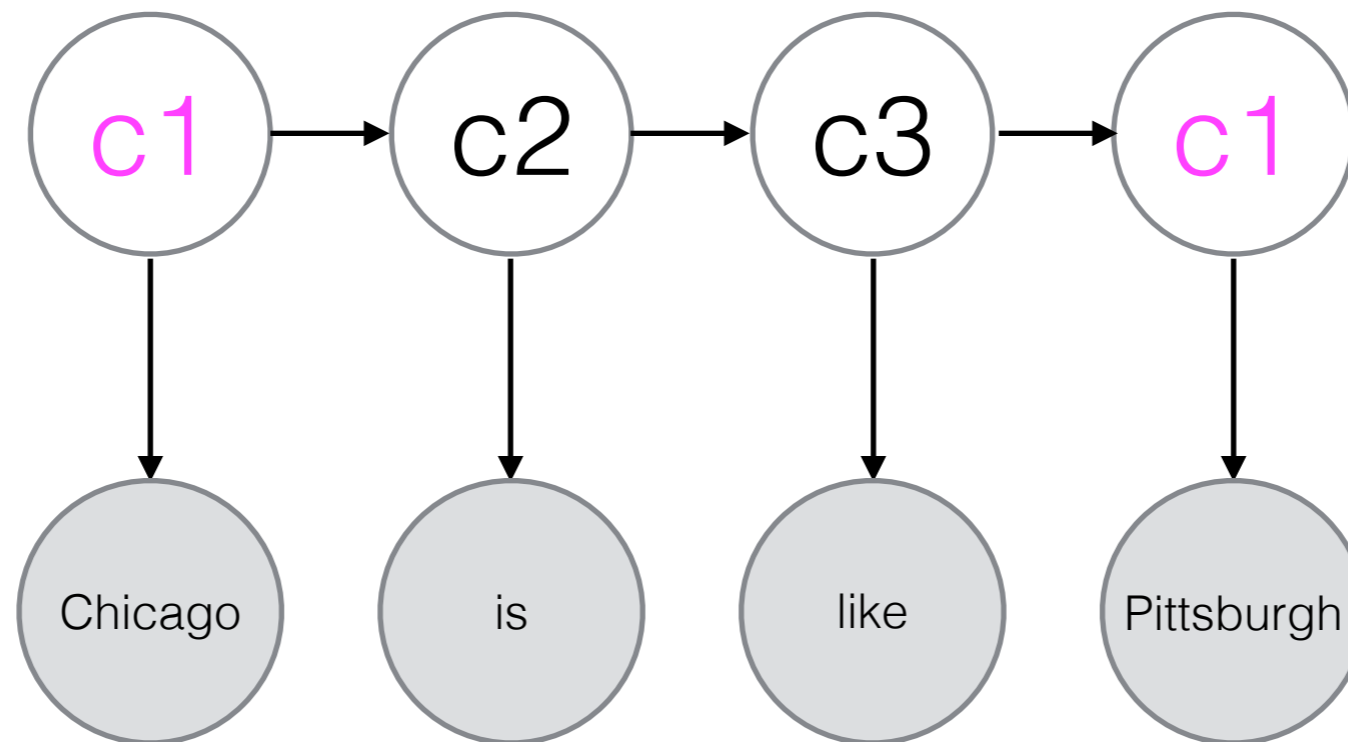
Representation Learning

“You shall know a word by the company it keeps” (Firth 1957)

- my boy's wicked smart
- my boy's hella smart
- my boy's very smart
- my boy's extremely smart
- my boy's _____ smart

Brown clustering

Unsupervised HMM, where each word type belongs to **one** class.



Brown clustering

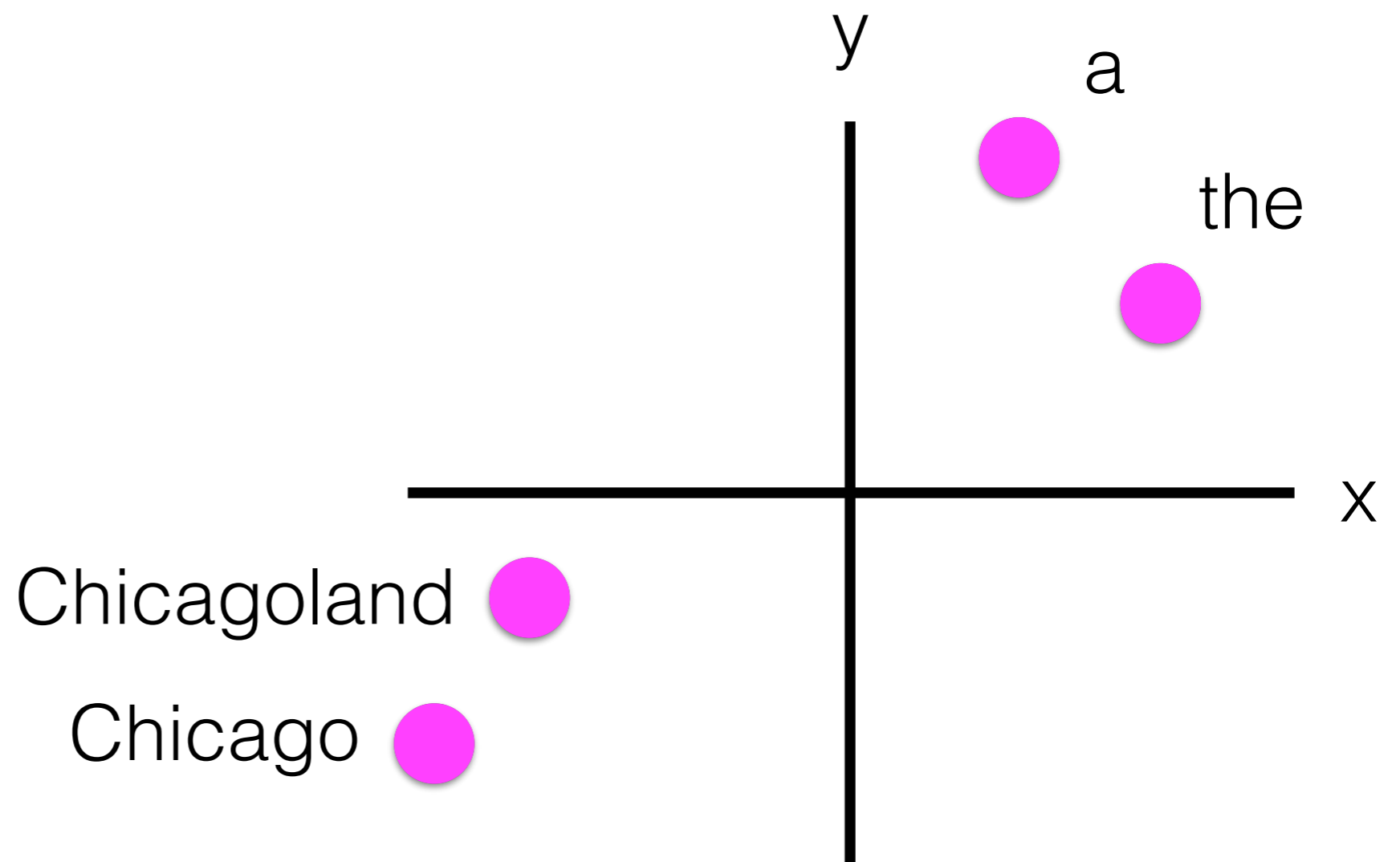
- Demo: 1000 clusters learned from 56M tweets
- http://www.ark.cs.cmu.edu/TweetNLP/cluster_viewer.html
- Code: <https://github.com/percyliang/brown-cluster>

Embeddings

- Represent each word in your vocabulary as a vector of K numbers

	x	y
the	2.1	2.5
a	1.5	3.7
Chicago	-3.0	-3.4
Chicagoland	-2.6	-0.5

Embeddings



Embeddings

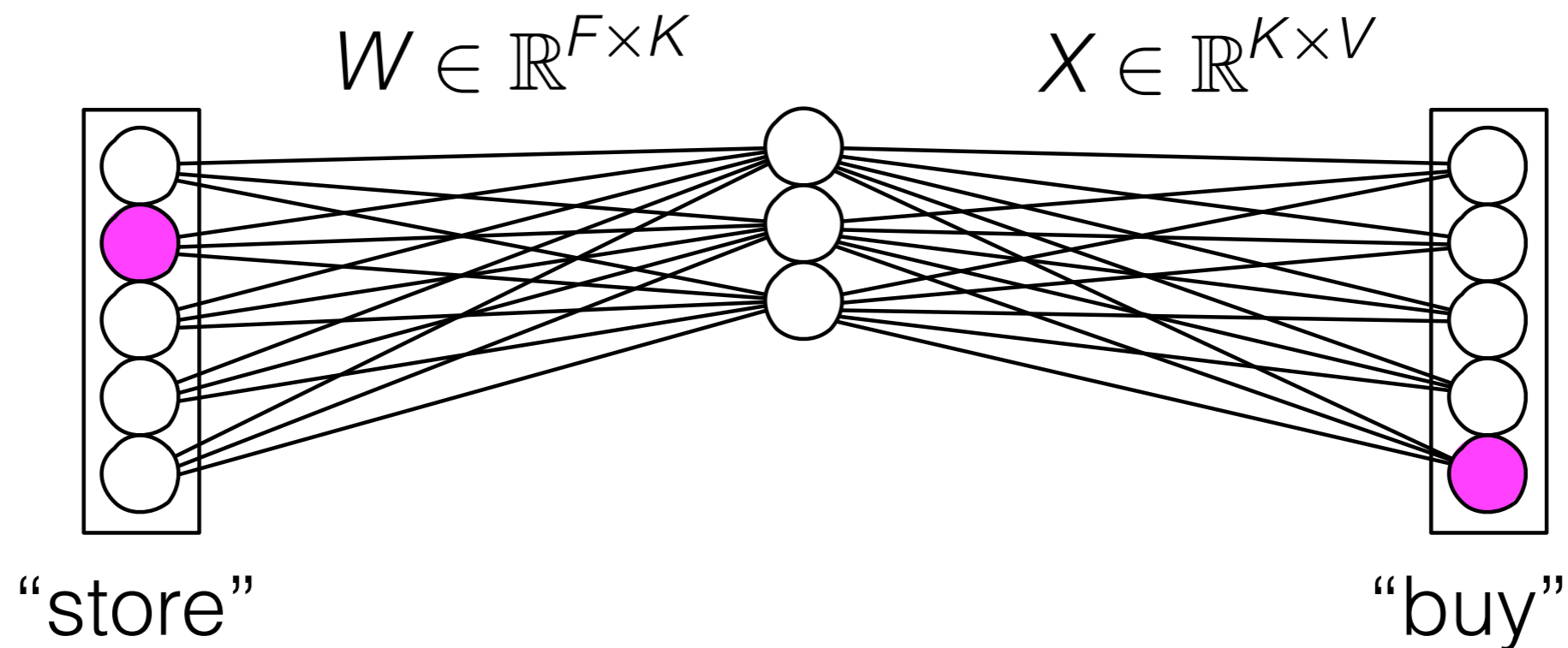
- Basic intuition: use a K-dimensional embedding for a word in a sentence to predict all of the words around it; find the value of the embedding to maximize your predictive accuracy.

Let's go to the ^{3.1}1.7 to buy some eggs.

Skip-Gram Embeddings

$$h_i = x^\top W_i$$

$$P(\text{word} = w|x, h, X) = \frac{\exp(h^\top X_w)}{\sum_v \exp(h^\top X_v)}$$



Embeddings

- Demo: <http://radimrehurek.com/2014/02/word2vec-tutorial/#app>
- Code: <https://code.google.com/p/word2vec/>

Word Representations

What do you do with word representations?

Word Representations

brown:169	brown:170	brown:171
Mr.	Chicago	New York
Mrs.	Chicagoland	NYC
	Chitown	NY

Word Representations

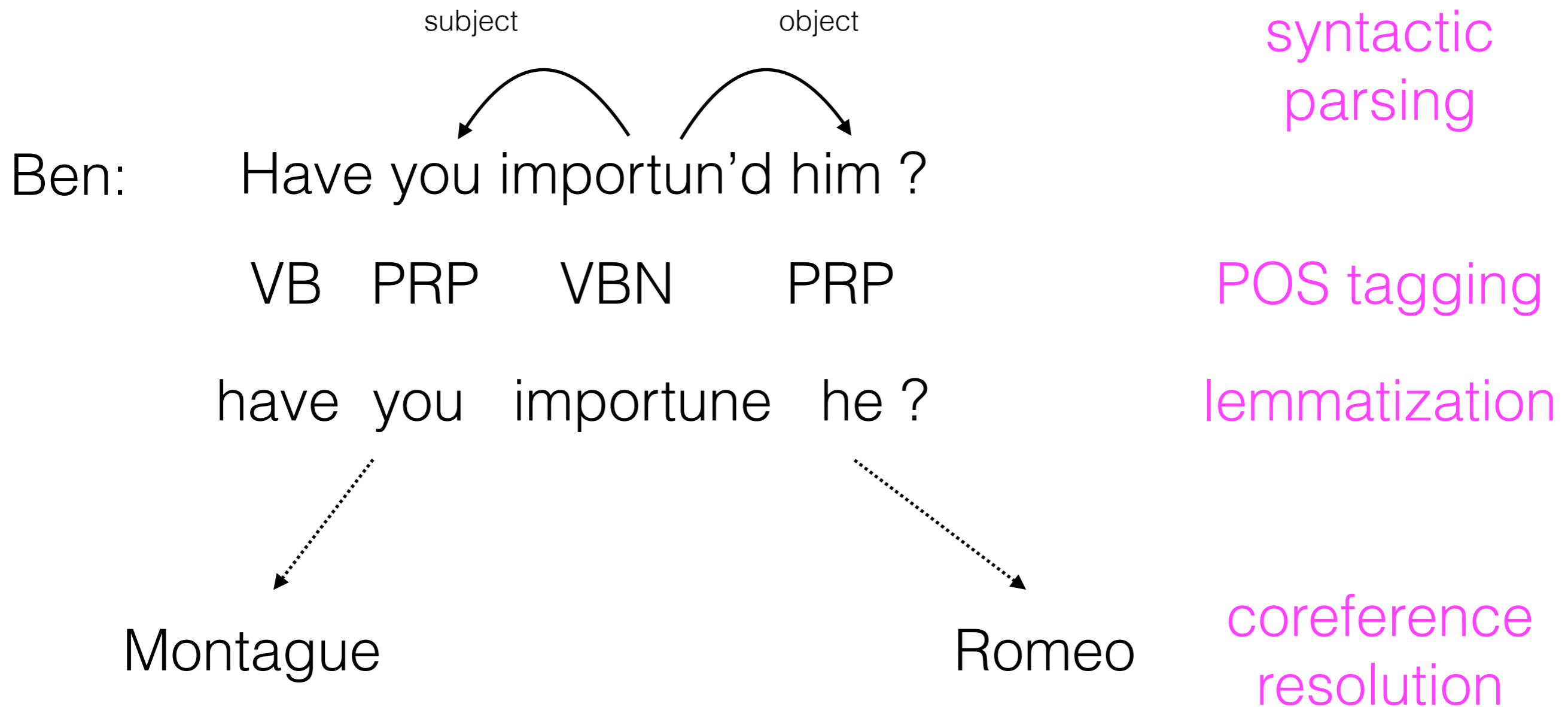
“I live in Chicago”

i	feat	value
1	I	1
2	live	1
3	in	1
4	New York	0
5	Chicago	1
6	Boston	0
7	Pittsburgh	0
8	snow	0
9	brown:170	1

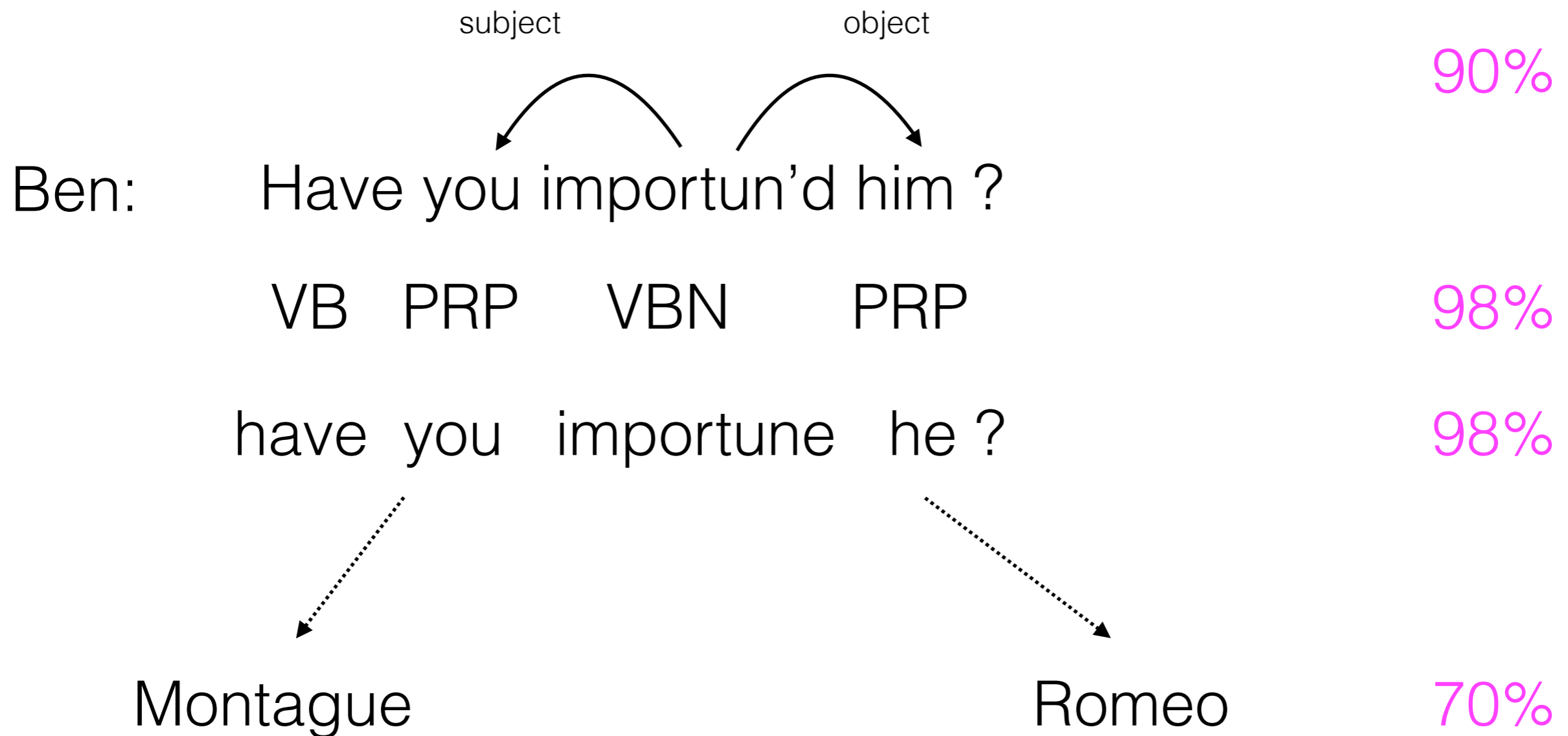
“I live in Chicagoland”

i	feat	value
1	I	1
2	live	1
3	in	1
4	New York	0
5	Chicago	0
6	Boston	0
7	Pittsburgh	0
8	snow	0
9	brown:170	1

NLP and beyond



NLP and beyond



NLP toolkits

- Tokenization, part of speech tagging, syntactic parsing, named entity recognition, coreference resolution.
- CoreNLP
<http://nlp.stanford.edu/software/corenlp.shtml>
- BookNLP
<https://github.com/dbamman/book-nlp>
- NLTK
<http://www.nltk.org>

Thanks!

- David Bamman
dbamman@cs.cmu.edu