

# Latent Dirichlet Allocation

Computational Linguistics

Alexander Koller

15 January 2019

with help from Christoph Teichmann  
and illustrations by Martín Villalba



# Today

- Today's lecture is about a method called *Latent Dirichlet Allocation (LDA)*.
- We care about it for two reasons:
  - ▶ It's an unsupervised method for identifying *topics* and words that are representative of them.
  - ▶ It's a showcase for a family of statistical models called *Bayesian models* which have many uses in CL.

# Let's start simple

- You and I are playing a coin-tossing game.  
I see you throw 63x H, 37x T.  
Should I believe that the coin is fair?
- Our model of the coin has one parameter,  $p = P(H)$ .
- Maximum-likelihood estimate:  $p = 0.63$ , i.e. not fair.
- But what about
  - ▶ my uncertainty about  $p$ ?
  - ▶ my prior beliefs about the fairness of the coin?

# Bayesian Models

- ML estimation and similar methods deliver *point estimates*: a single value for each parameter that optimizes some criterion.
  - ▶ Likelihood:  $P(\text{observations} \mid \text{parameters})$
- Bayesian models: estimate a *probability distribution*  $P(\text{parameters} \mid \text{observations})$  over parameters.
  - ▶ assume a *prior* over parameters, which encodes beliefs in parameter values before making any observations
  - ▶ update prior to *posterior* after making some observations
  - ▶ uncertainty about parameter values is reflected at all times in the pd

# The Dirichlet distribution

- Take the parameter  $p$  itself as the value of a random variable.
  - ▶ need a probability distribution over real numbers;  
more specifically, over tuples of numbers that sum to one
- We use the *Dirichlet distribution*.

$p_1, \dots, p_K \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$  means:

$$P(p_1, \dots, p_K) = \frac{1}{B(\alpha)} (p_1^{\alpha_1-1} \cdot \dots \cdot p_K^{\alpha_K-1})$$

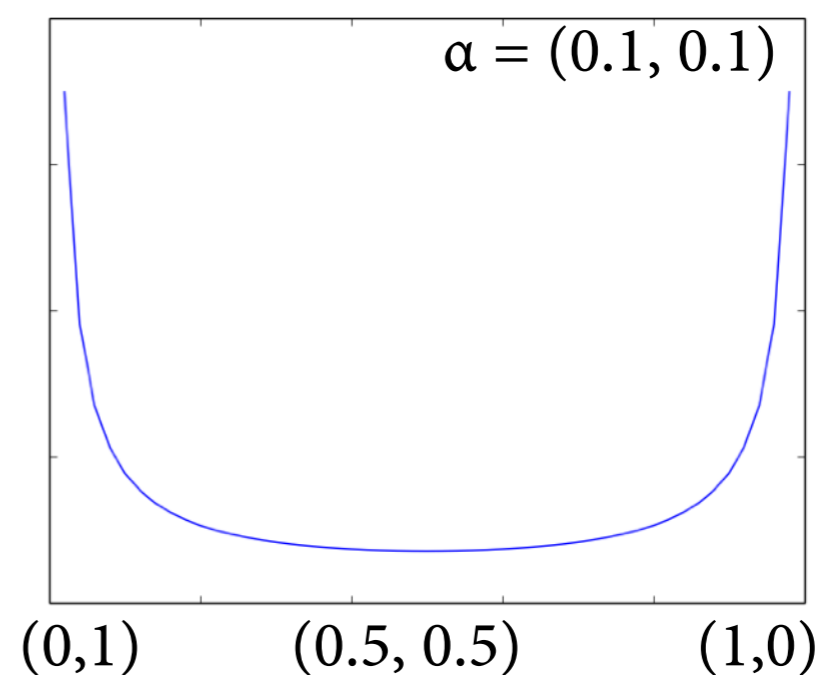
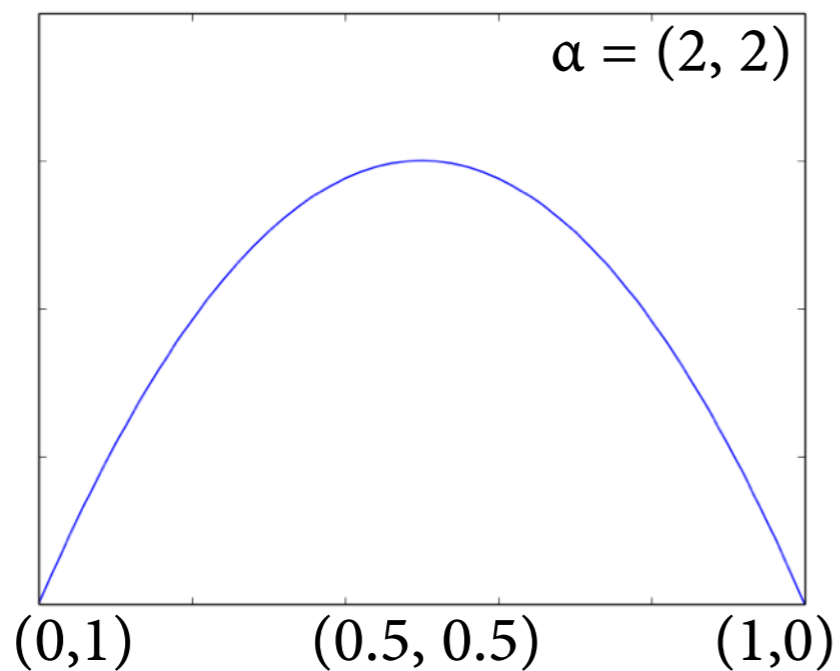
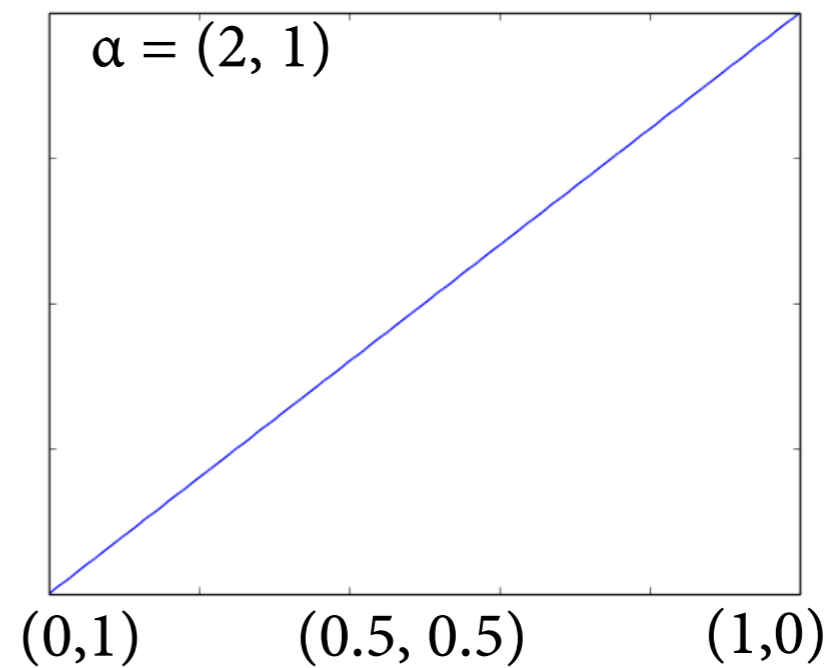
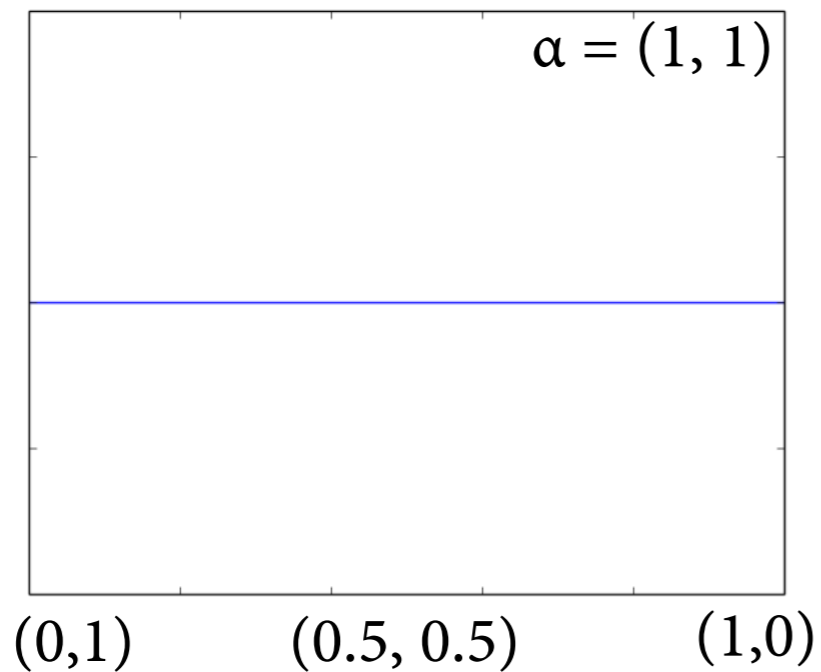
Dir only defined if  
the  $p_i$  sum to 1

this is the *beta function*  
(needed to normalize to 1)

$\alpha_1, \dots, \alpha_K$  are called  
*hyperparameters*

# Dirichlet distributions, $K = 2$

$$P(p_1, \dots, p_K) = \frac{1}{B(\alpha)} (p_1^{\alpha_1 - 1} \cdot \dots \cdot p_K^{\alpha_K - 1})$$



# Bayesian parameter estimation

- We are interested in  $P(M)$  over our model  $M = (p)$ . This model is very simple; will make more complex later.
- Before we make any observations, we have a *prior distribution*:  $P(M) = \text{Dir}_{\alpha, \alpha}(p, 1-p)$
- We can then *update* this to a *posterior distribution* based on observed data:

$$P(M | D) = \frac{P(D | M) \cdot P(M)}{P(D)} \propto P(D | M) \cdot P(M)$$

posterior                      likelihood                      prior

# Calculating posteriors

prior:  $P(p) = \text{Dir}_{\alpha, \alpha}(p, 1 - p) \propto p^{\alpha-1} \cdot (1 - p)^{\alpha-1}$

likelihood:  $P(i \times \text{H}, k \times \text{T} \mid p) = p^i \cdot (1 - p)^k$

posterior: 
$$\begin{aligned} P(p \mid i \times \text{H}, k \times \text{T}) &\propto P(i \times \text{H}, k \times \text{T} \mid p) \cdot P(p) \\ &\propto p^i \cdot (1 - p)^k \cdot p^{\alpha-1} \cdot (1 - p)^{\alpha-1} \\ &= p^{i+\alpha-1} \cdot (1 - p)^{k+\alpha-1} \end{aligned}$$

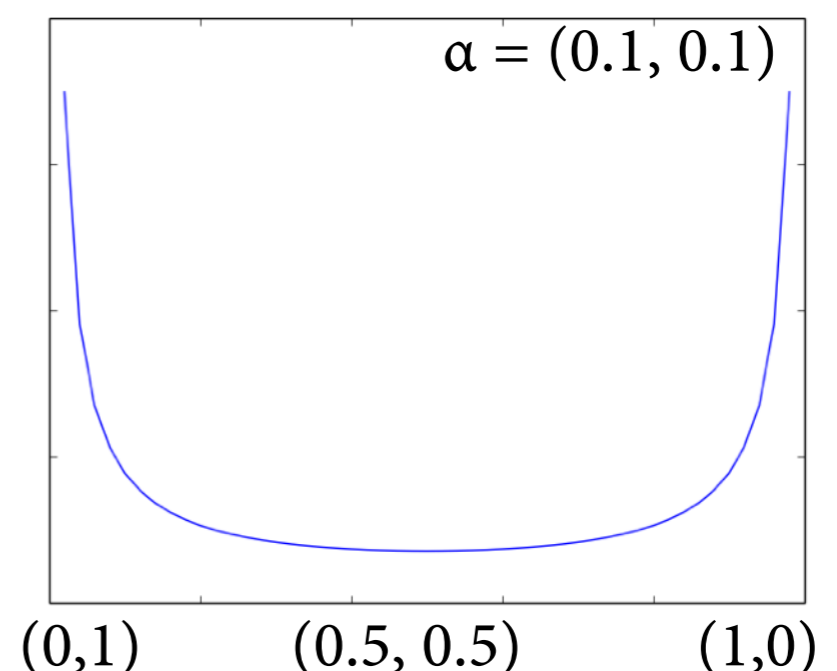
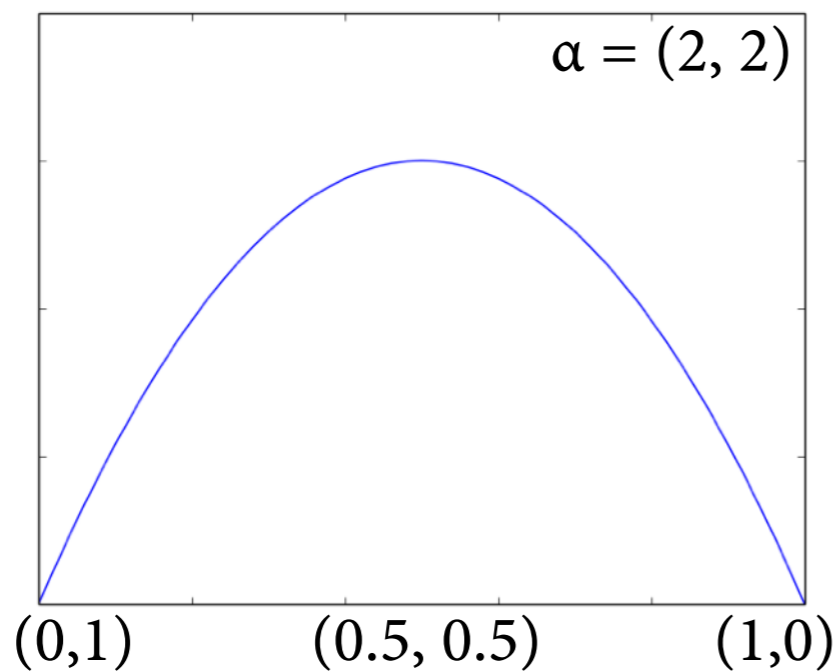
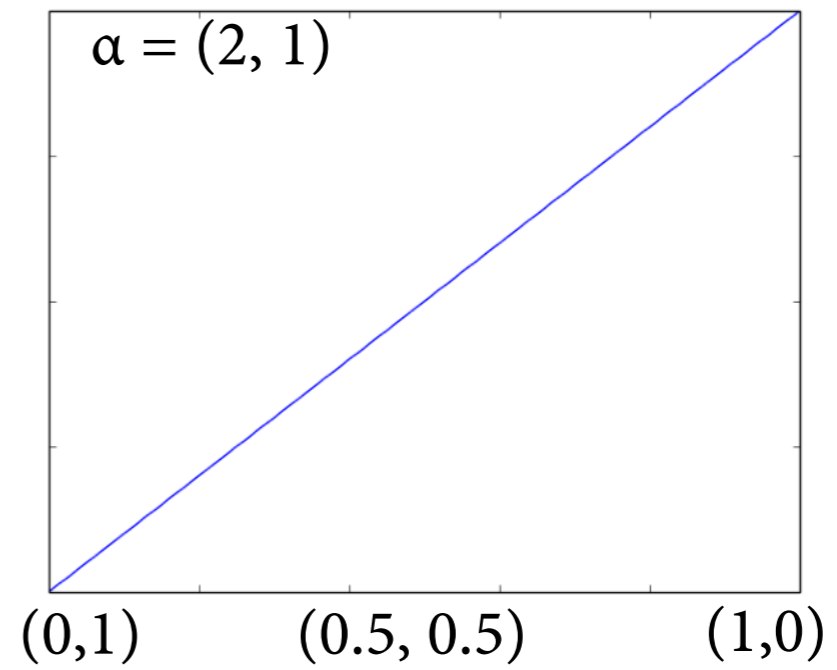
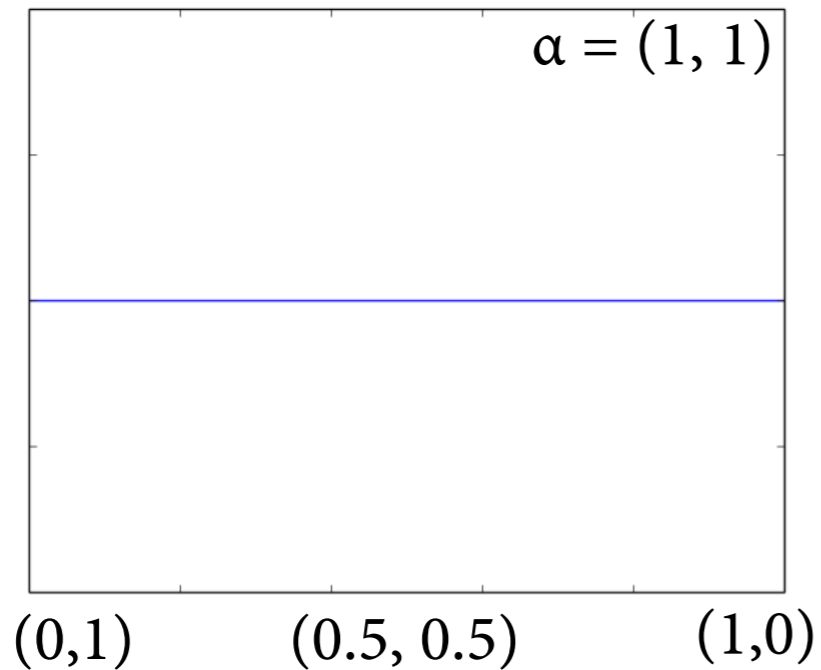
More precisely, we have:

$$P(p \mid i \times \text{H}, k \times \text{T}) = \text{Dir}_{\alpha+i, \alpha+k}(p, 1 - p)$$



# The Dirichlet distribution

$$P(p_1, \dots, p_K) = \frac{1}{B(\alpha)} (p_1^{\alpha_1 - 1} \cdot \dots \cdot p_K^{\alpha_K - 1})$$

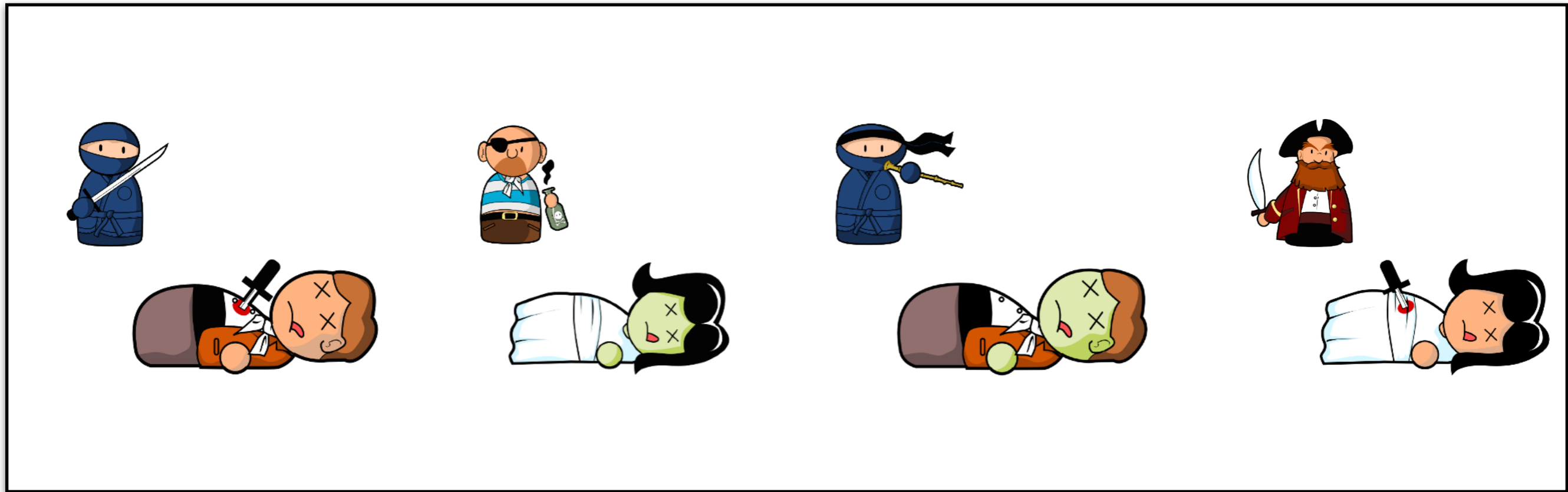


# Conjugate distributions

- Crucially,  $P(M)$  and  $P(M | D)$  have the same shape (product of Dirichlets). This is because Dirichlet and Categorical are *conjugate distributions*.
  - ▶ because  $K = 2$  for the coin, we really only used the Beta (not Dirichlet) and Bernoulli (not Categorical) distributions
- This is makes the math very convenient.
- The hyperparameters of the Dirichlets are updated by adding the observed counts to the hp. of the priors.
  - ▶ priors thus perform smoothing in a very principled way

# The next step

Say you come across some people who have been stabbed or poisoned.  
You know that each of them was killed by a pirate or a ninja.  
You can tell how each person died, but not by whom they were killed.



# Our task

- We observe  $N$  people with their causes of death.
- Questions we are interested in:
  - ▶ Who killed each villager?  
 $z_1, \dots, z_N \in \{\text{pi}, \text{ni}\}$
  - ▶ How many were killed by pirates, how many by ninjas?  
 $P(\text{pi}) = \theta_{\text{pi}}, P(\text{ni}) = \theta_{\text{ni}}; \text{ thus, } \theta_{\text{pi}} + \theta_{\text{ni}} = 1$
  - ▶ How likely is it that a pirate chooses to stab someone?  
 $P(\text{st} \mid \text{pi}) = \phi_{\text{st}|\text{pi}}; \text{ thus, } P(\text{po} \mid \text{pi}) = \phi_{\text{po}|\text{pi}} = 1 - \phi_{\text{st}|\text{pi}}$
  - ▶ How likely is it that a ninja chooses to stab someone?  
 $P(\text{st} \mid \text{ni}) = \phi_{\text{st}|\text{ni}}; \text{ thus, } P(\text{po} \mid \text{ni}) = \phi_{\text{po}|\text{ni}} = 1 - \phi_{\text{st}|\text{ni}}$

# Fundamental approach

- Goal: Bayesian model with parameters  $\theta$ ,  $\phi_{\text{pi}}$ ,  $\phi_{\text{ni}}$ .
  - ▶ maximum likelihood: try to estimate concrete values for each parameter
  - ▶ Bayesian: estimate *probability distribution*  $P(\theta, \phi_{\text{pi}}, \phi_{\text{ni}})$
- In practice, the model will have *latent variables*  $z$ , which cannot be observed directly (e.g. pirate/ninja).
- Will marginalize over model parameters and work with  $P(z \mid \text{observations})$  directly.

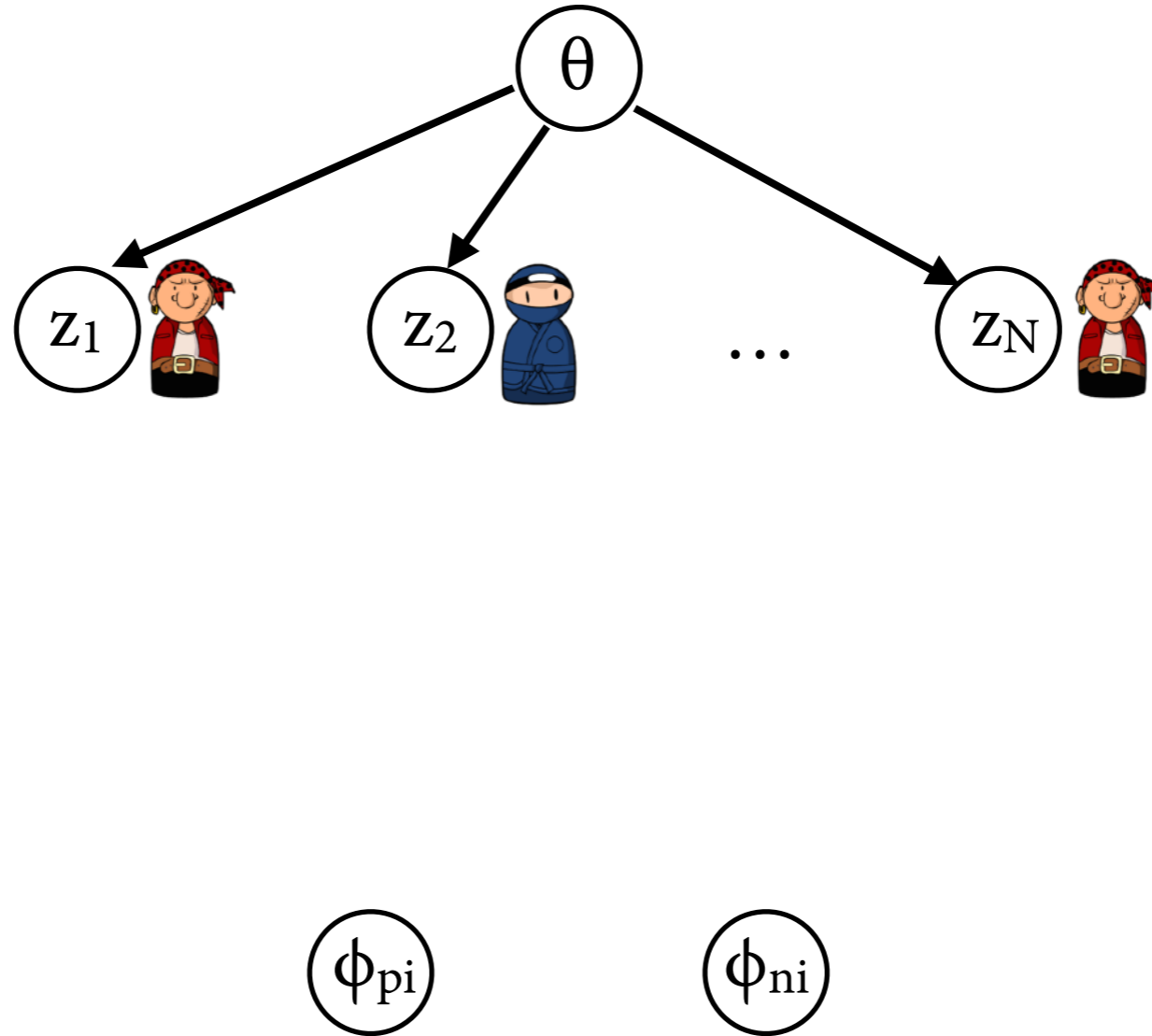
# Generative story: Idea

$\theta$

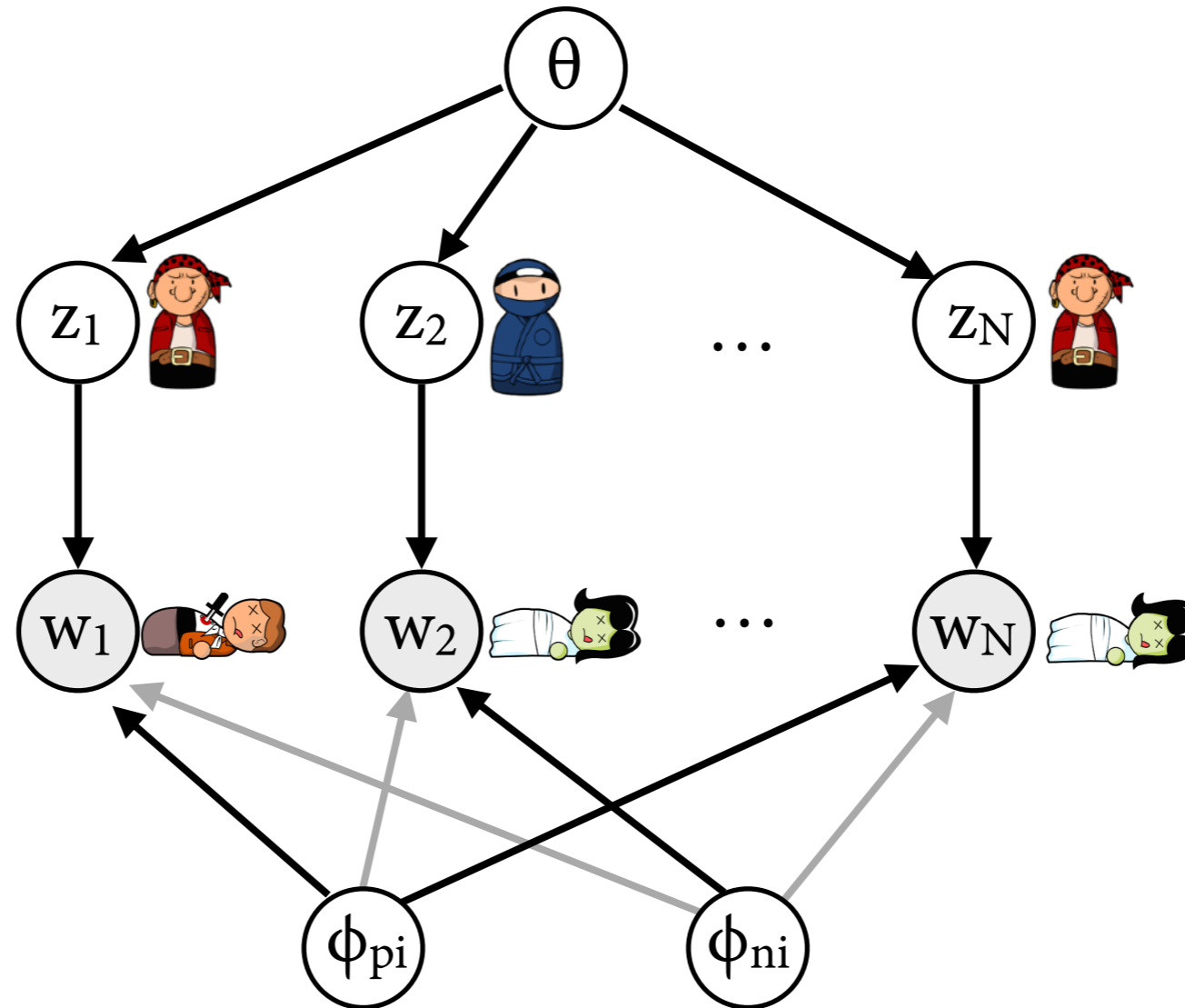
$\phi_{pi}$

$\phi_{ni}$

# Generative story: Idea

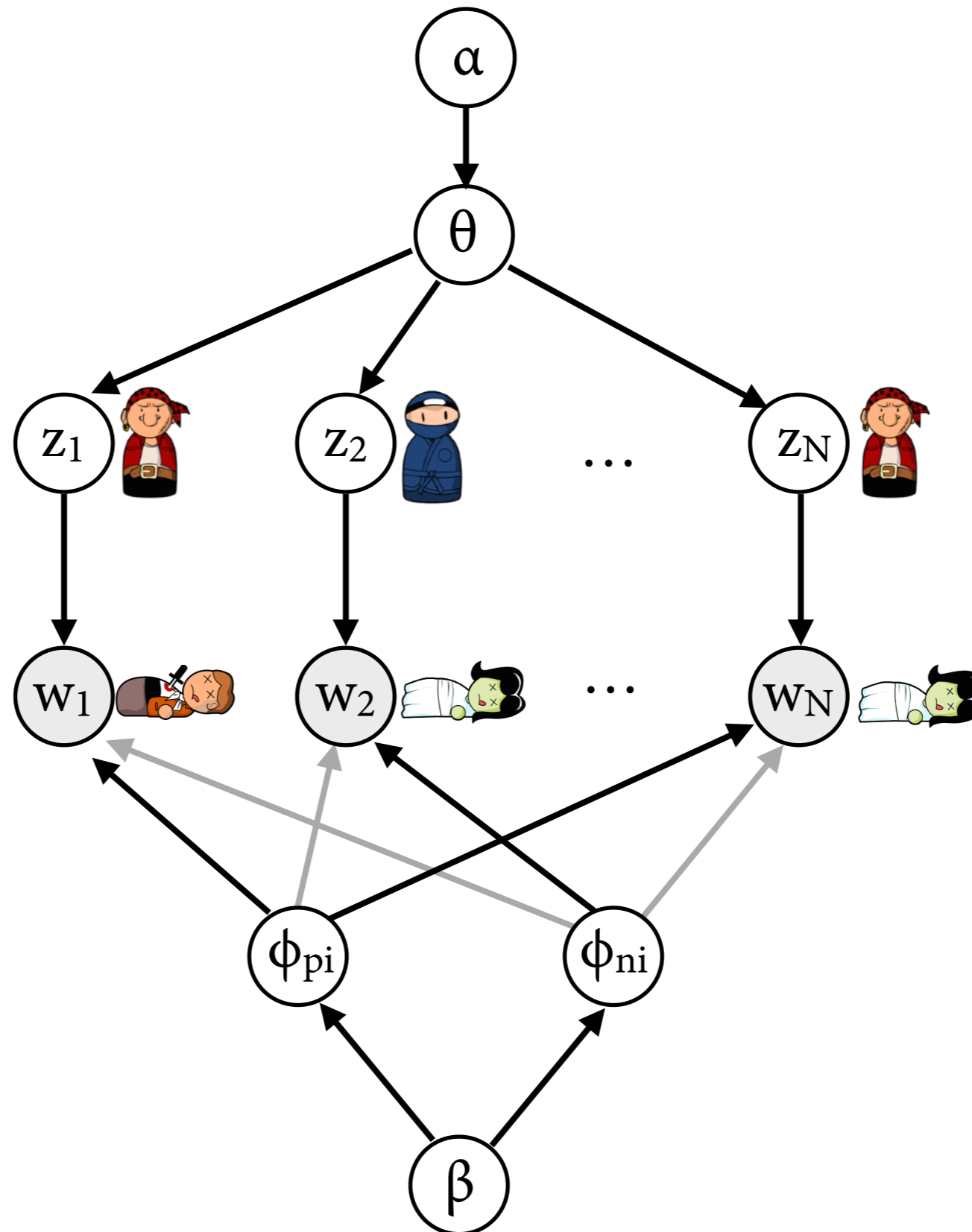


# Generative story: Idea





# Generative story: Idea



# Generative story

- We assume deaths are generated as follows:

$$(\theta_{pi}, \theta_{ni}) \sim \text{Dir}(\alpha, \alpha)$$

$$(\phi_{st|pi}, \phi_{po|pi}), (\phi_{st|ni}, \phi_{po|ni}) \sim \text{Dir}(\beta, \beta)$$

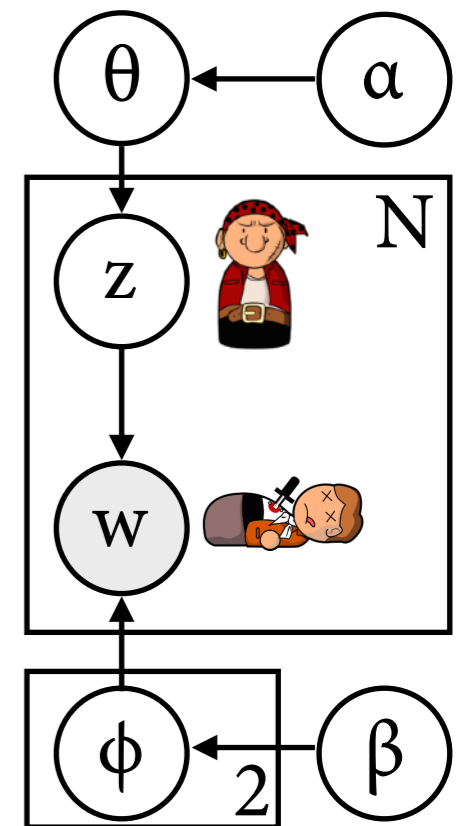
$$z_1, \dots, z_K \sim \text{Categorical}(\theta)$$

$$w_i \sim \text{Categorical}(\phi_{z_i})$$

- That is:

- ▶  $P(z_i = pi) = \theta_{pi}$ ,  $P(z_i = ni) = \theta_{ni}$


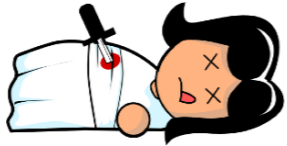


- ▶ if  $z_i$  came out as “pi”, then  $P(w_i = st) = \phi_{st|pi}$



I abbreviate  $\theta = (\theta_{pi}, \theta_{ni})$ ,  $\phi_{pi} = (\phi_{st|pi}, \phi_{po|pi})$ ,  $\phi_{ni} = (\phi_{st|ni}, \phi_{po|ni})$ .  
 $\alpha, \beta$  are assumed given and are called *hyperparameters*.

# Supervised learning

If all killers are known,  $P(M | D)$  is easy to compute.

$i$	$z_i$	$w_i$
1		
2		

$$P(M) = \text{Dir}_{\alpha, \alpha}(\theta) \cdot \text{Dir}_{\beta, \beta}(\phi_{pi}) \cdot \text{Dir}_{\beta, \beta}(\phi_{ni})$$

$$\propto \theta_{pi}^{\alpha-1} \cdot \theta_{ni}^{\alpha-1} \cdot \phi_{st|pi}^{\beta-1} \cdot \phi_{po|pi}^{\beta-1} \cdot \phi_{st|ni}^{\beta-1} \cdot \phi_{po|ni}^{\beta-1}$$

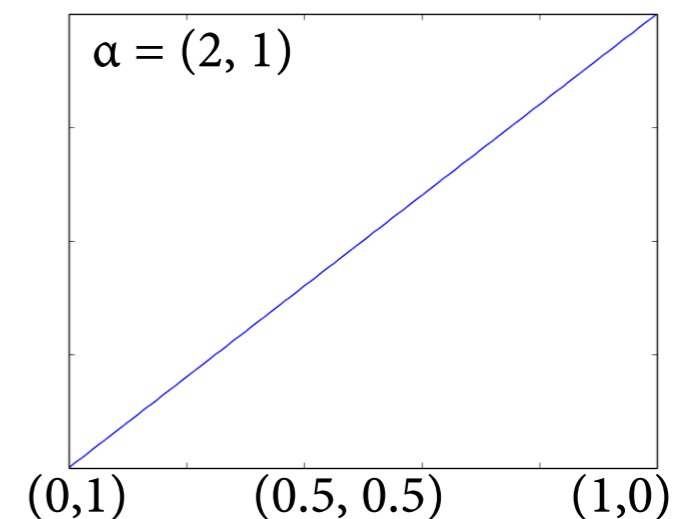
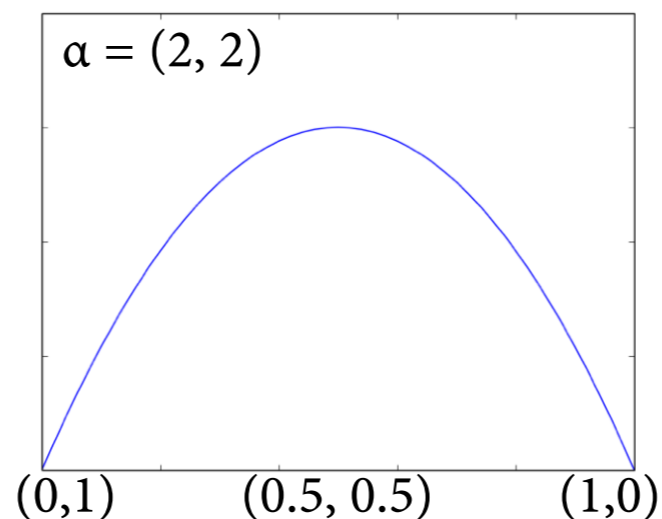
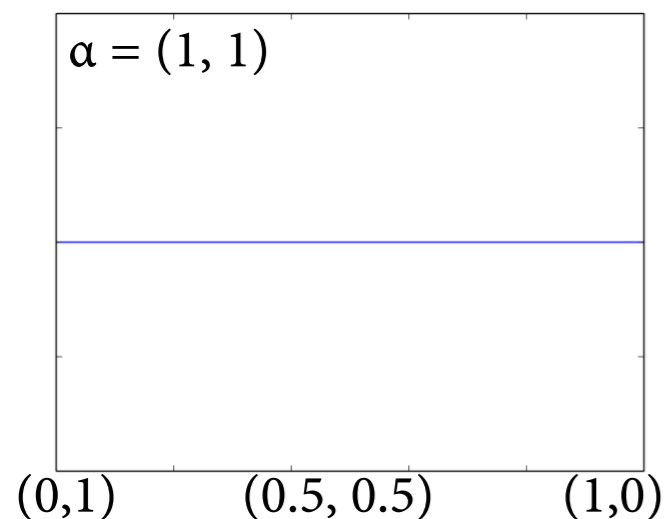
$$P(D | M) = P(z_1 = pi, w_1 = st, z_2 = ni, w_2 = po)$$

$$= \theta_{pi} \cdot \phi_{st|pi} \cdot \theta_{ni} \cdot \phi_{po|ni}$$

$$P(M | D) \propto P(D | M) \cdot P(M)$$

$$\propto \theta_{pi}^{\alpha} \cdot \theta_{ni}^{\alpha} \cdot \phi_{st|pi}^{\beta} \cdot \phi_{po|pi}^{\beta-1} \cdot \phi_{st|ni}^{\beta-1} \cdot \phi_{po|ni}^{\beta}$$



$$\propto \text{Dir}_{\alpha+1, \alpha+1}(\theta) \cdot \text{Dir}_{\beta+1, \beta}(\phi_{pi}) \cdot \text{Dir}_{\beta, \beta+1}(\phi_{ni})$$



# Unsupervised learning

- In the original scenario, we can only observe deaths, not killers. Then  $P(D | M)$  is less convenient:

$$\begin{aligned} P(D | M) &= P(w_1 = \text{st}, w_2 = \text{po} | M) \\ &= \sum_{k_1, k_2 \in \{\text{pi}, \text{ni}\}} P(z_1 = k_1, w_1 = \text{st}, z_2 = k_2, w_2 = \text{po} | M) \end{aligned}$$

i	$z_i$	$w_i$
1	??	
2	??	

- This sums over a number of terms that is exponential in  $N$ , and thus infeasible to compute.
- $M = (\theta, \phi_{\text{pi}}, \phi_{\text{ni}})$

# Latent variables

- Many interesting quantities can be expressed in terms of distribution over the latent variables.

$$P(z | w) = \int P(z, M | w) dM = \int P(z | M, w) \cdot P(M | w) dM$$

- Some examples:

ninja/pirate mixing proportion

$$\frac{1}{N} \cdot E_{P(z|w)}[C(z_i = \text{ninja})]$$

pirate habits

$$E_{P(z|w)}[C(z_i = \text{pirate}, w_i = \text{stab})] / E_{P(z|w)}[C(z_i = \text{pirate})]$$

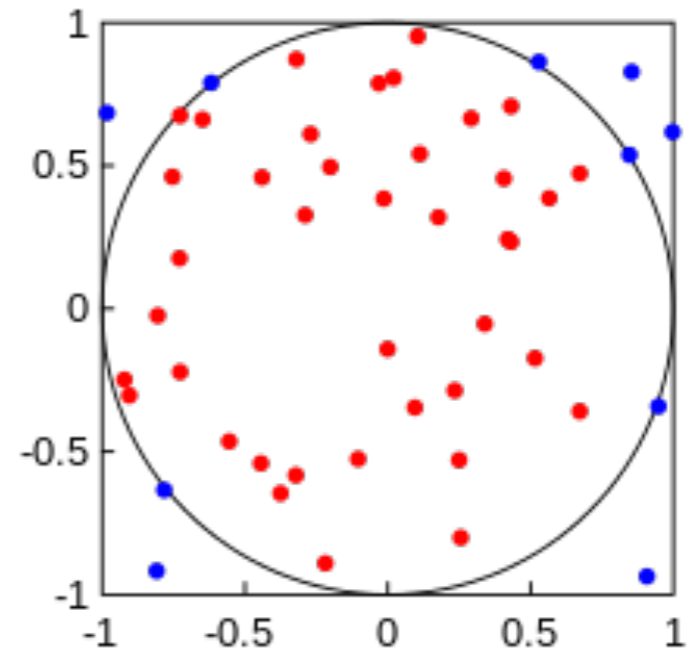
probability that first villager was killed by a pirate

$$E_{P(z|w)}[ \|z_1 = \text{pirate}\| ]$$

# Estimating expected values

- Expected values can be approximated by *sampling*.  
To compute  $E_{P(X)}[f(X)]$ :
  - ▶ draw  $S$  samples  $x^{(1)}, \dots, x^{(S)}$  from  $P(X)$
  - ▶ estimate  $E[f(X)] \approx \frac{1}{S} \cdot \sum_{i=1}^S f(x^{(i)})$
- Example: To estimate  $\pi$ , sample points from square and count how many fall into the circle.

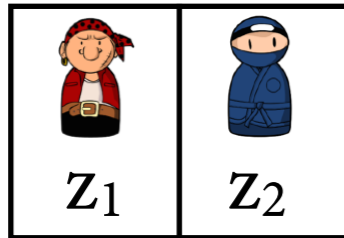
$$\pi/4 \approx E_{P(x,y)}[ \|x^2 + y^2 \leq 1\| ]$$



# EVs under latent variables

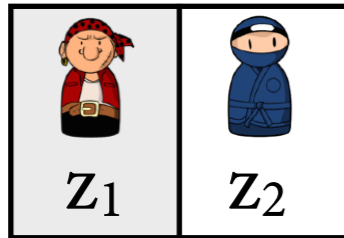
- We could estimate expected values under  $P(z | w)$  using sampling. However,  $P(z | w)$  is usually of a form that makes direct sampling difficult.
- Instead, we can use *Gibbs sampling*:
  - ▶ Start from an initial guess  $z_1, \dots, z_N$  for the latent variables.
  - ▶ Repeatedly resample guess for some  $z_i$  conditioned on all other  $z$ 's, i.e. from  $P(z_i | w, z_{-i})$ . This is much easier than sampling from  $P(z | w)$  itself.
  - ▶ Can prove that probability of observing a sample for  $z$  as a whole converges to  $P(z | w)$ .

# Gibbs Sampling

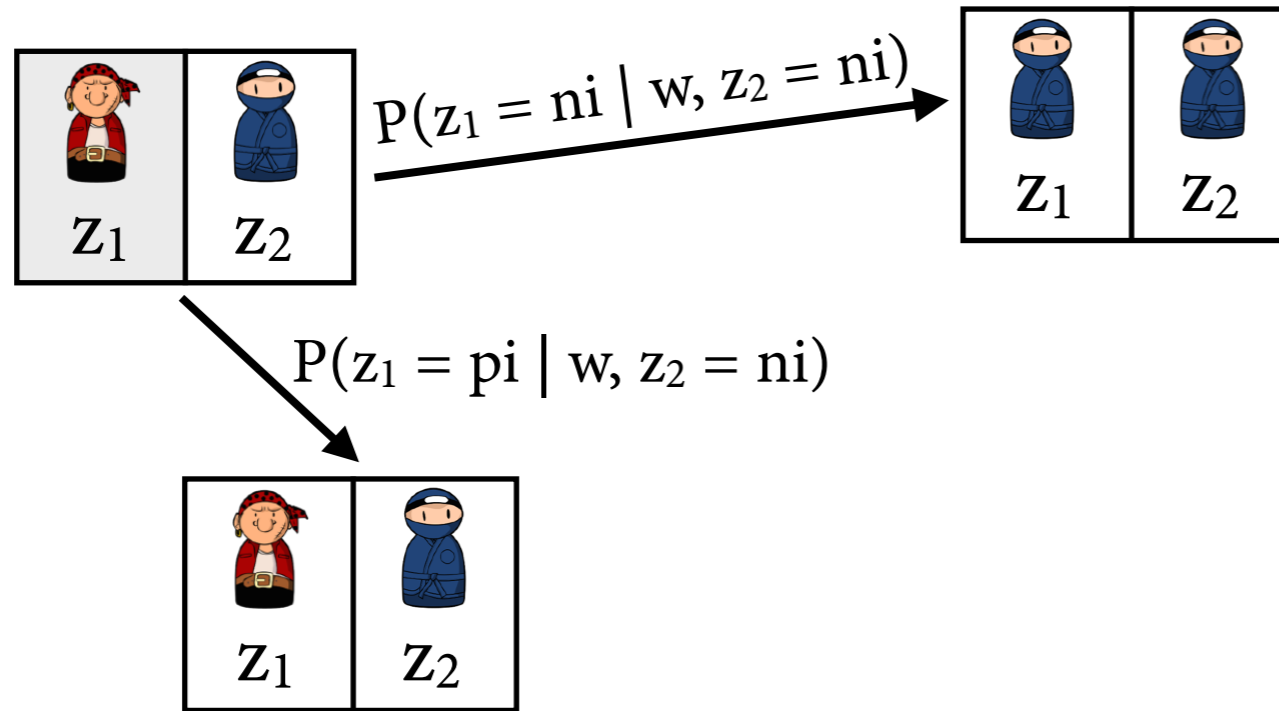




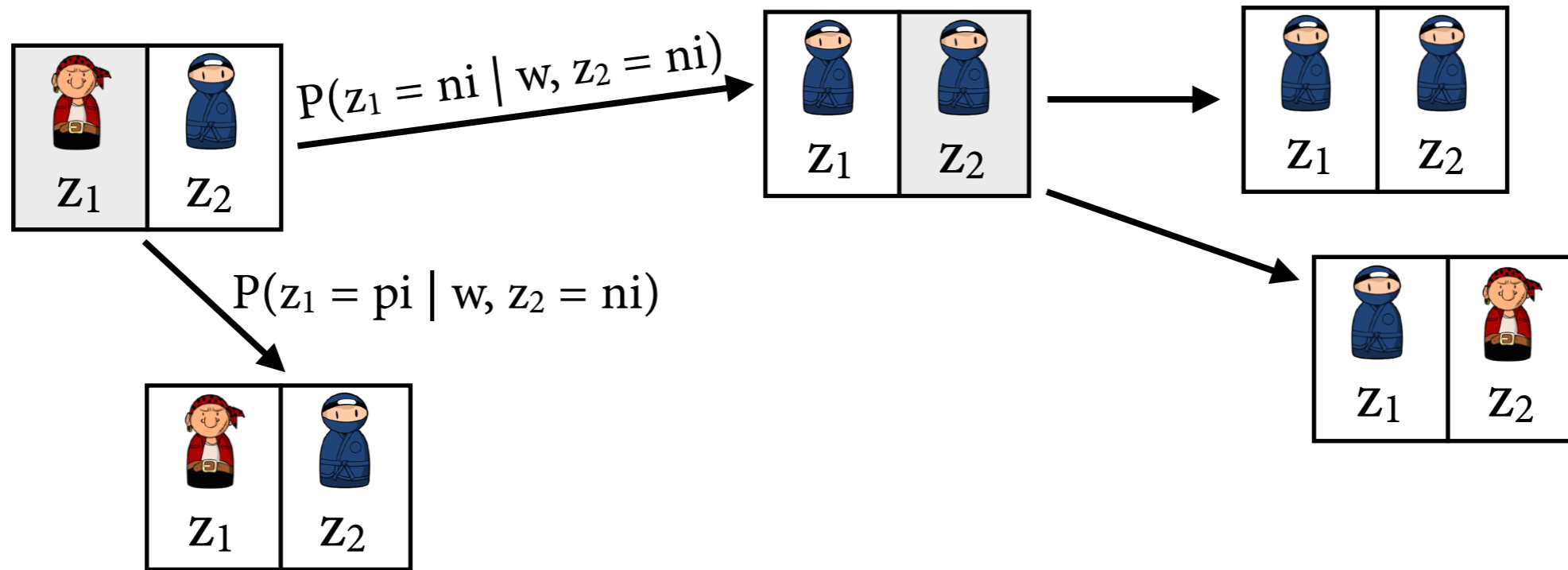
# Gibbs Sampling



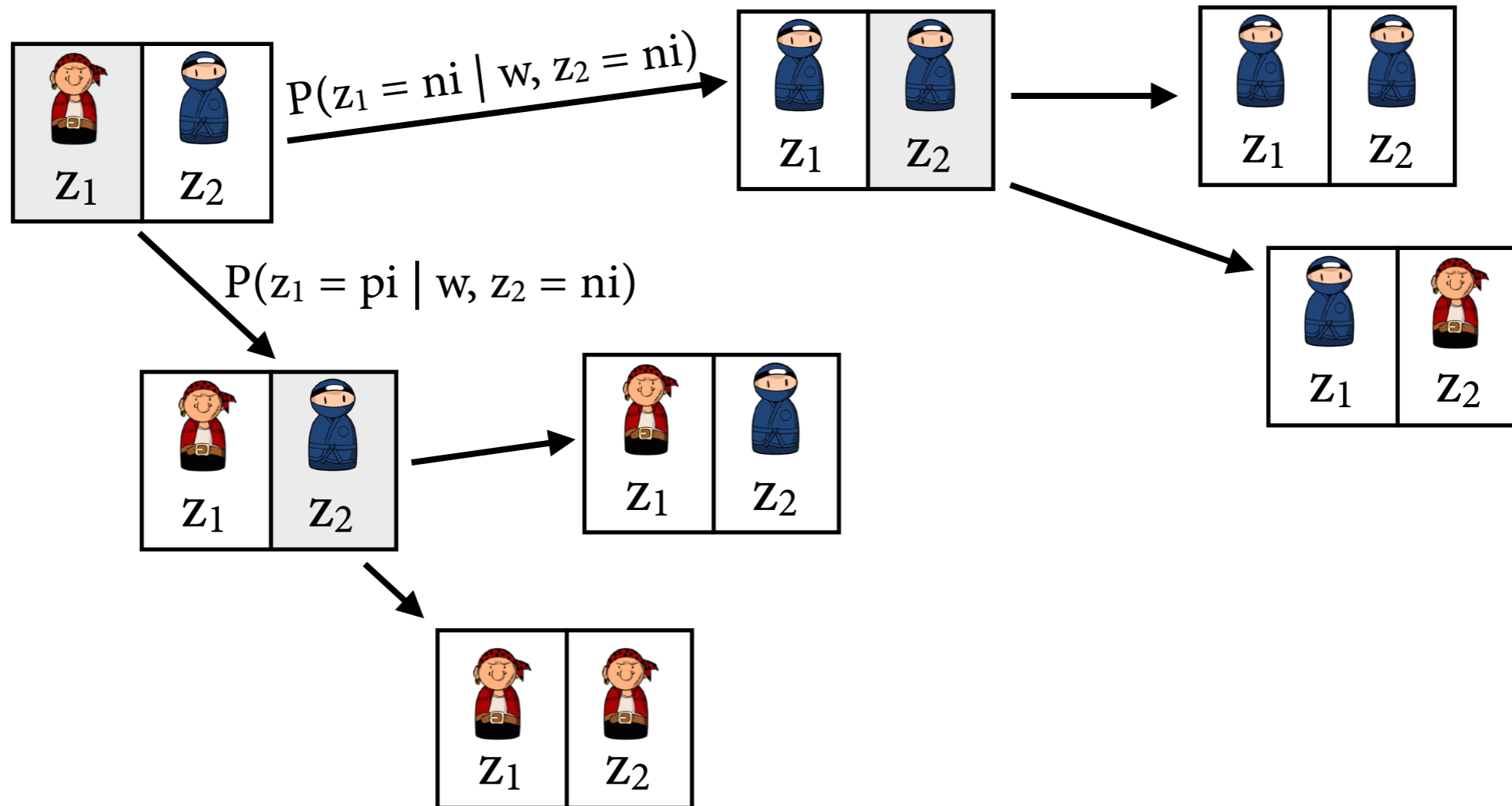
# Gibbs Sampling



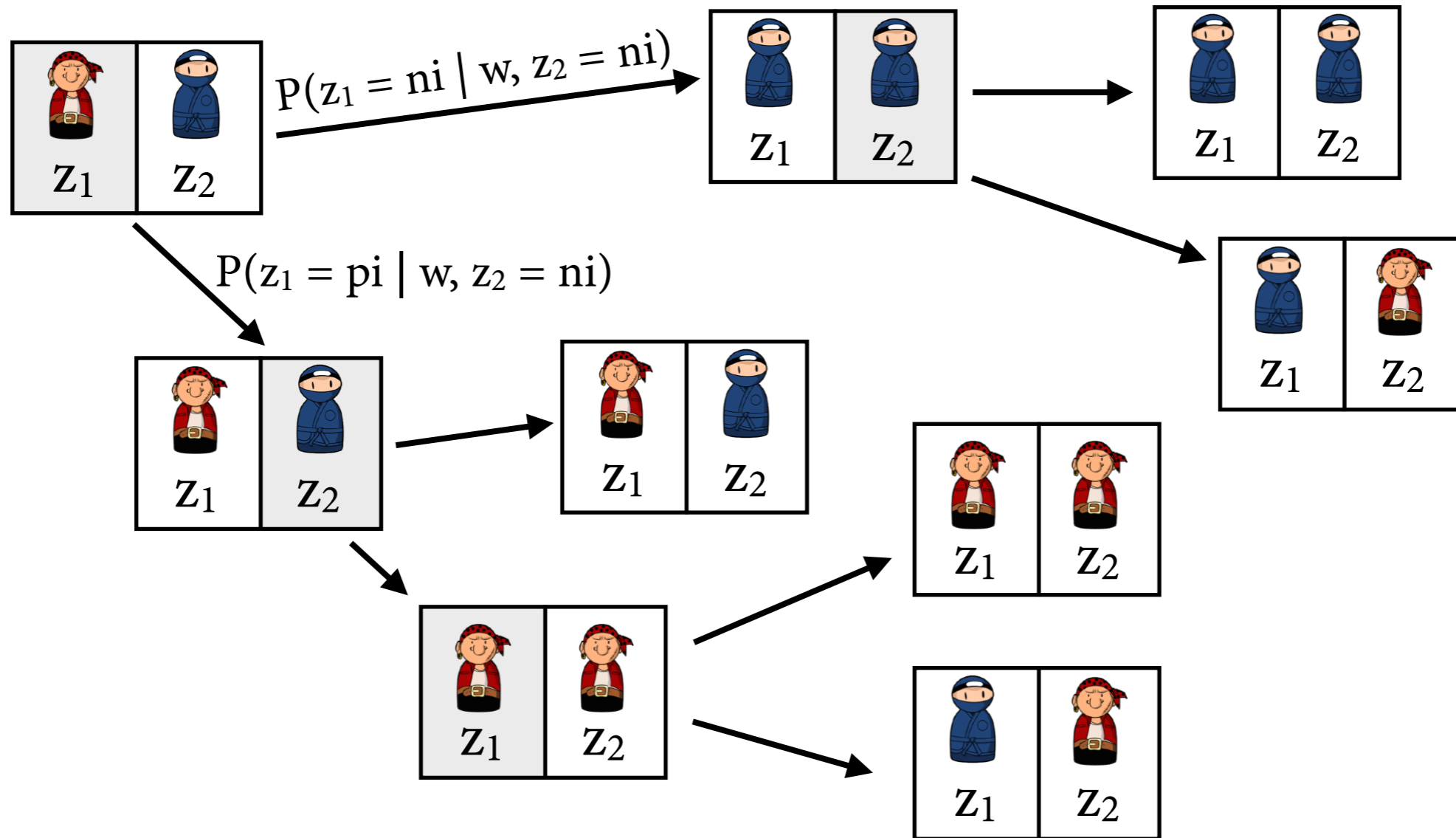
# Gibbs Sampling



# Gibbs Sampling



# Gibbs Sampling



# Transition probabilities

- It remains to determine the transition probabilities  $P(z_i \mid w, z_{-i})$ .
- Formula turns out to be remarkably simple:

$$\begin{aligned} P(z_i = p_i \mid w, z_{-i}) &\propto P(w, z_{-i}, z_i = p_i) \\ &= \int \int P(w, z_{-i}, z_i = p_i, \theta, \phi) \, d\theta \, d\phi \\ &= \dots \end{aligned}$$

$$\propto (n_{p_i}^{(-i)} + \alpha_{p_i}) \frac{n_{p_i, w_i}^{(-i)} + \beta_{w_i | p_i}}{\sum_{w'} n_{p_i, w'}^{(-i)} + \beta_{w' | p_i}}$$

↖  
# people other than i that  
were killed by pirates  
in current sample

↙  
# people other than i  
that were killed by pirates  
using method w'

# Topic models

### Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

### Documents

#### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK— How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

### Topic proportions and assignments

(Blei, Comm. ACM 12)

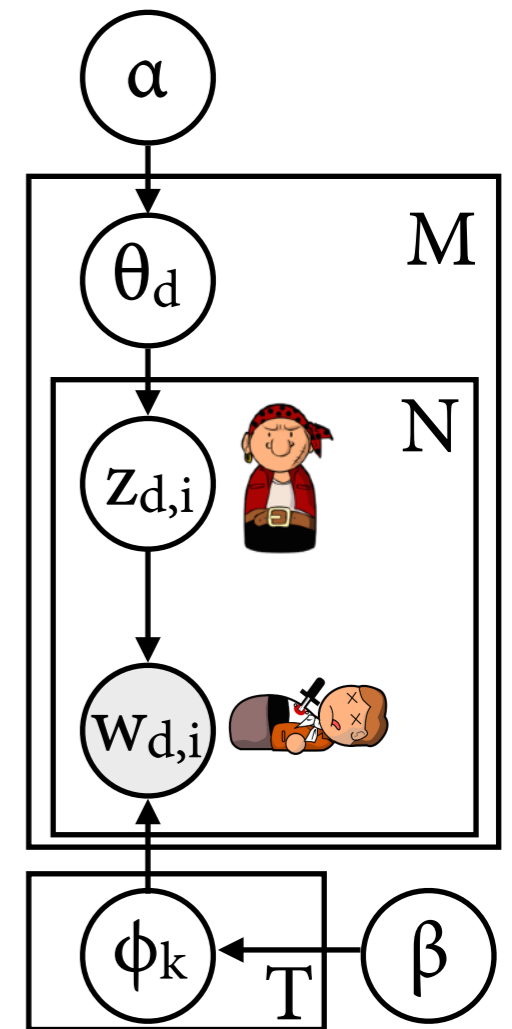
learn: word probs.  
for (abstract) topics

← given: raw documents →

learn: topic mixture  
in each document

# Latent Dirichlet Allocation

- Topic modeling is almost the same problem as the pirate/ninja problem:
  - ▶ abstract topics = {pirate, ninja}
  - ▶ words in document = {stabbed, poisoned}
- Full LDA makes two changes:
  - ▶ can have  $T$  topics instead of just two, and also more than two different words
  - ▶ there are  $M > 1$  *documents*, and each document can have its own mixture  $\theta_d$  of topics





# Gibbs sampler for LDA

prob of reassigning  
token #i as topic t

# t occurs with word  $w_i$   
except at position i

# t occurs in document  
that contains position i,  
except at position i

$$P(z_i = t \mid z_{-i}, w) \propto \frac{n_{-i,t}^{(w_i)} + \beta}{n_{-i,t}^{(\cdot)} + W \cdot \beta} \cdot \frac{n_{-i,t}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + T \cdot \alpha}$$

# t occurs anywhere in corpus,  
except at position i

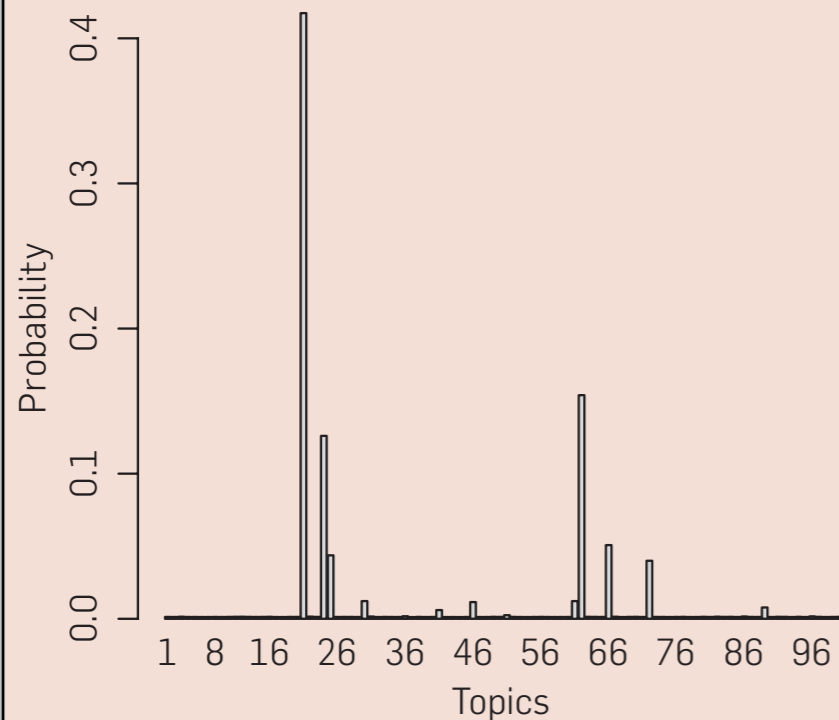
# tokens in that document,  
minus one (for position i)

$W$  = vocabulary size /  $T$  = number of topics

(Griffiths & Steyvers 2004)

# Examples

(Blei 2012)



## “Genetics”

human  
genome  
dna  
genetic  
genes  
sequence  
gene  
molecular  
sequencing  
map  
information  
genetics  
mapping  
project  
sequences

## “Evolution”

evolution  
evolutionary  
species  
organisms  
life  
origin  
biology  
groups  
phylogenetic  
living  
diversity  
group  
new  
two  
common

## “Disease”

disease  
host  
bacteria  
diseases  
resistance  
bacterial  
new  
strains  
control  
infectious  
malaria  
parasite  
parasites  
united  
tuberculosis

## “Computers”

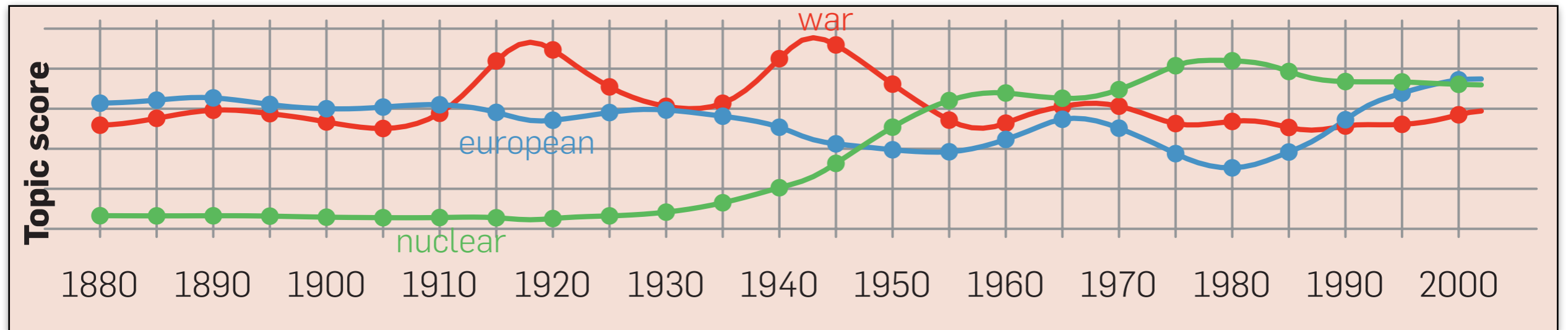
computer  
models  
information  
data  
computers  
system  
network  
systems  
model  
parallel  
methods  
networks  
software  
new  
simulations

topic mixture for  
one article in *Science*

15 words with highest  $\phi_{k,w}$   
for each topic over whole corpus  
(with made-up topic label)

# Examples

development of topics from *Science* over time (1880-2002)



# Conclusion

- LDA and extensions for topic modeling.
  - ▶ Topics interesting in their own right, also useful in various applications.
  - ▶ Simplest useful Bayesian model in NLP.
- We used (collapsed) Gibbs sampling to approximate expected values.
  - ▶ Alternative is *Variational Bayes*: approximate  $P(M|D)$  on paper, then solve integral exactly.
- Limitation: Number  $T$  of topics must be given. We will fix this next time.