

Universität Potsdam
Institut für Informatik
Lehrstuhl Maschinelles Lernen



Graphical Models

Niels Landwehr

Agenda

- Graphical models: syntax and semantics.
- Inference in graphical models (exact, approximate)
- Graphical models in machine learning
 - ◆ Recap: Coin Tosses
 - ◆ Recap: Bayesian linear regression
 - ◆ Latent Dirichlet allocation
 - ◆ Hidden Markov models

Recap: Bayesian Linear Regression

- Solving regression problems

$$L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

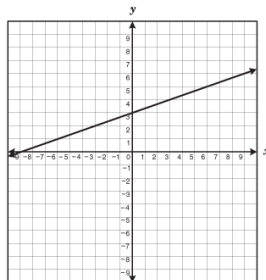
$$\mathbf{X} = \begin{pmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nm} \end{pmatrix} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$\mathbf{x}_i \in \mathbb{R}^m$ feature vector

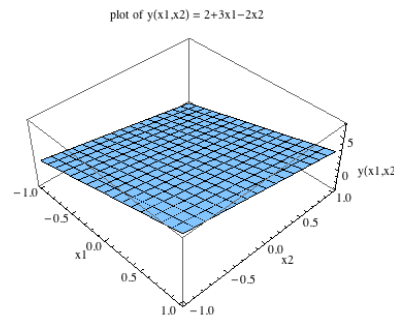
$y_i \in \mathbb{R}$ real-valued target

- Linear regression

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{i=1}^m w_i x_i$$



\mathbf{w} „parameter vector“, „weight vector“



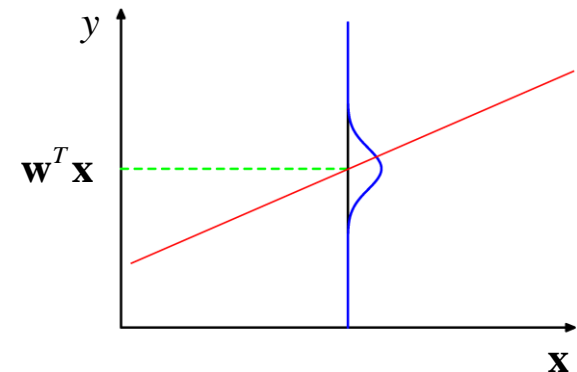
Recap: Bayesian Linear Regression

- Discriminative setting: \mathbf{x}_i fixed input, y_i generated from \mathbf{x}_i and \mathbf{w} plus Gaussian noise.

$$\begin{aligned} p(y | \mathbf{x}, \mathbf{w}) &= \mathbf{w}^T \mathbf{x} + N(y | 0, \sigma^2) \\ &= N(y | \mathbf{w}^T \mathbf{x}, \sigma^2) \end{aligned}$$

$$y_i \sim p(y | \mathbf{x}_i, \mathbf{w})$$

discriminative model: $p(\mathbf{x})$ not modeled



- Bayesian approach: posterior \propto prior \times likelihood

$$\underbrace{p(\mathbf{w} | L)}_{\text{posterior}} \propto \underbrace{p(L | \mathbf{w})}_{\text{likelihood}} \underbrace{p(\mathbf{w})}_{\text{prior}}$$

Recap: Bayesian Linear Regression

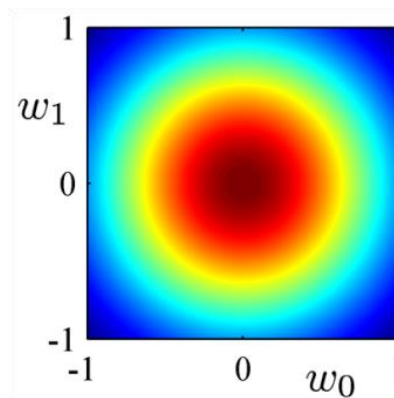
- Likelihood of data under model \mathbf{w} :

$$\begin{aligned} p(\mathbf{y} | \mathbf{X}, \mathbf{w}) &= \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2) \quad i.i.d. \\ &= \prod_{i=1}^N N(y_i | \mathbf{w}^T \mathbf{x}_i, \sigma^2) \end{aligned}$$

- Normally distributed prior over models:

$$p(\mathbf{w}) = N(\mathbf{w} | \mathbf{0}, \tau^2 I)$$

$$I = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$



Isotropic multivariate normal distribution, mean $\mathbf{0}$, variance τ^2

Bayesian Linear Regression as a Graphical Model

- What are the random variables?
 - ◆ Labels y_1, \dots, y_N , model \mathbf{w}
 - ◆ Not: $\mathbf{x}_1, \dots, \mathbf{x}_N$, hyperparameters σ^2, τ^2
 - ◆ Inputs \mathbf{x}_i behave like hyperparameters (fixed quantities)
- Joint distribution over labels and parameter vector

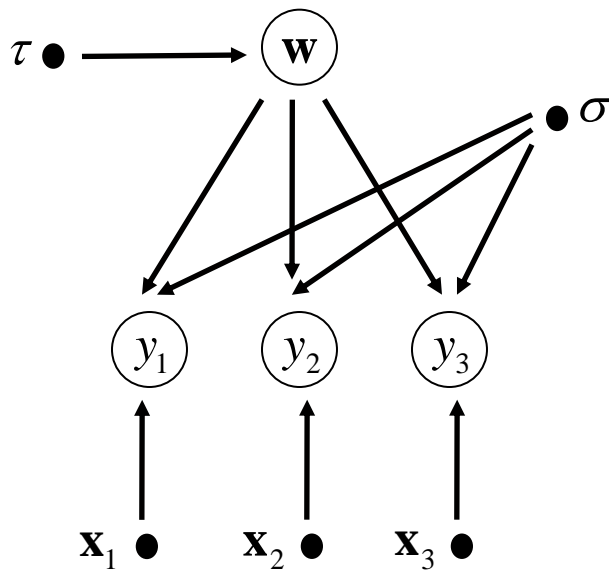
$$\begin{aligned} p(y_1, \dots, y_N, \mathbf{w} \mid \mathbf{x}_1, \dots, \mathbf{x}_N, \sigma^2, \tau^2) &= \overbrace{p(\mathbf{w} \mid \tau^2)}^{\text{Prior}} \overbrace{p(y_1, \dots, y_N \mid \mathbf{w}, \mathbf{x}_1, \dots, \mathbf{x}_N, \sigma^2)}^{\text{Likelihood}} \\ &= p(\mathbf{w} \mid \tau^2) \prod_{i=1}^N p(y_i \mid \mathbf{w}, \mathbf{x}_i, \sigma^2) \end{aligned}$$

- Representation of Bayesian linear regression as a graphical model: structure can be seen from the form of the joint distribution.

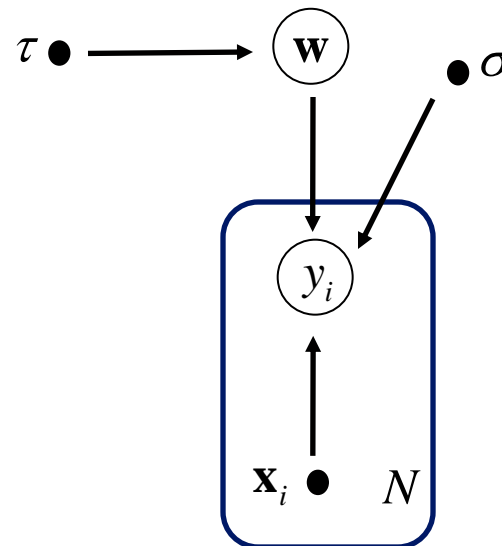
Bayesian Linear Regression as a Graphical Model

$$p(y_1, \dots, y_N, \mathbf{w} \mid \mathbf{x}_1, \dots, \mathbf{x}_N, \sigma^2, \tau^2) = p(\mathbf{w} \mid \tau^2) \prod_{i=1}^N p(y_i \mid \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

Graphical model, $N=3$



Graphical model, Plate notation



Bayes-optimal Prediction

- When applying model, need prediction for novel test instance \mathbf{x} :

$$\mathbf{x} \mapsto y$$

- Bayesian prediction

$$y_* = \arg \max_y p(y | \mathbf{x}, \mathbf{y}, \mathbf{X}, \sigma^2, \tau^2)$$

$$= \arg \max_y \int p(y | \mathbf{x}, \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \sigma^2, \tau^2) d\mathbf{w} \quad \text{„Bayesian model averaging“}$$

Bayesian Linear Regression as a Graphical Model

- Graphical model representation of Bayesian linear regression, including a novel test instance.

$$p(y_1, \dots, y_N, y, \mathbf{w} \mid \mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}, \sigma^2, \tau^2) = p(\mathbf{w} \mid \tau^2) \left(\prod_{i=1}^N p(y_i \mid \mathbf{w}, \mathbf{x}_i, \sigma^2) \right) p(y \mid \mathbf{w}, \mathbf{x}, \sigma^2)$$

Graphical model, $N=3$

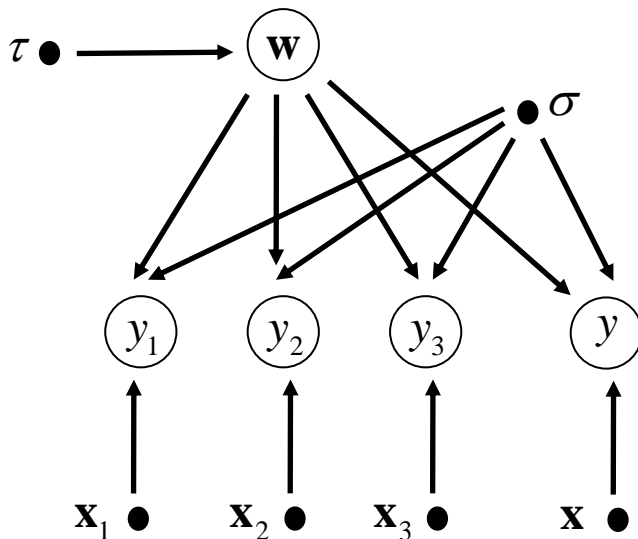
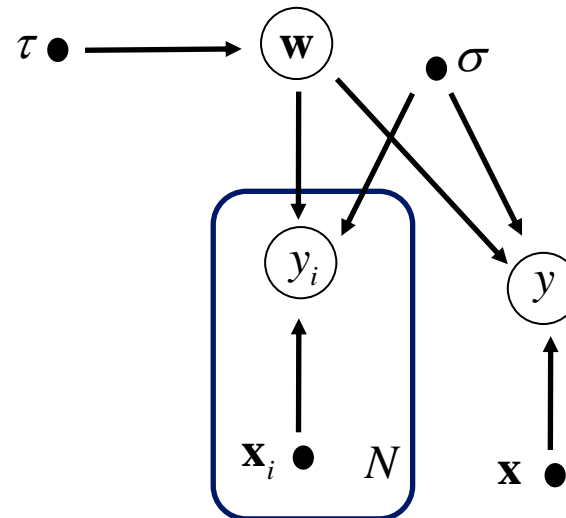
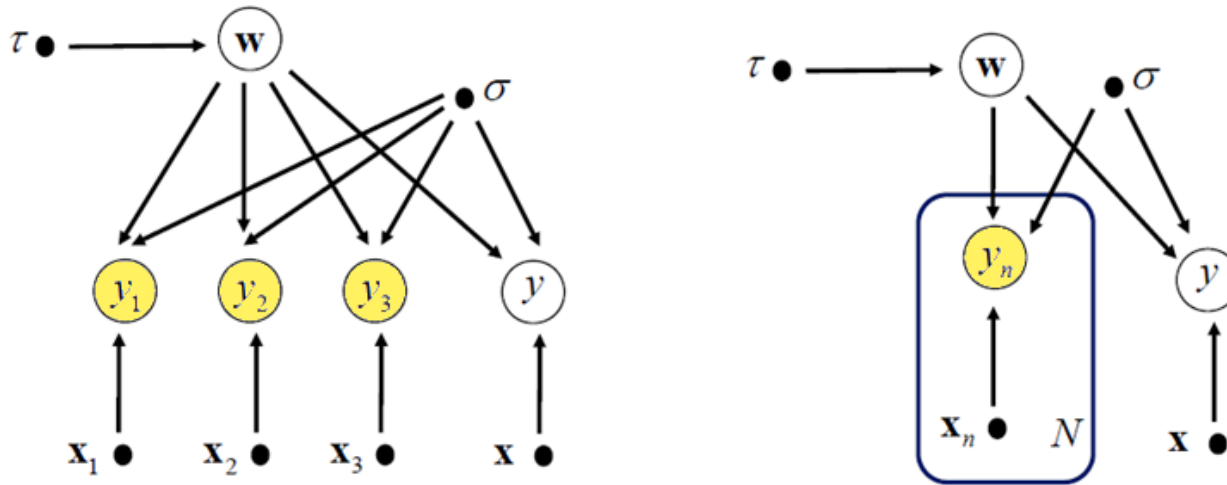


Plate notation



Bayesian Linear Regression as a Graphical Model



■ Bayesian prediction

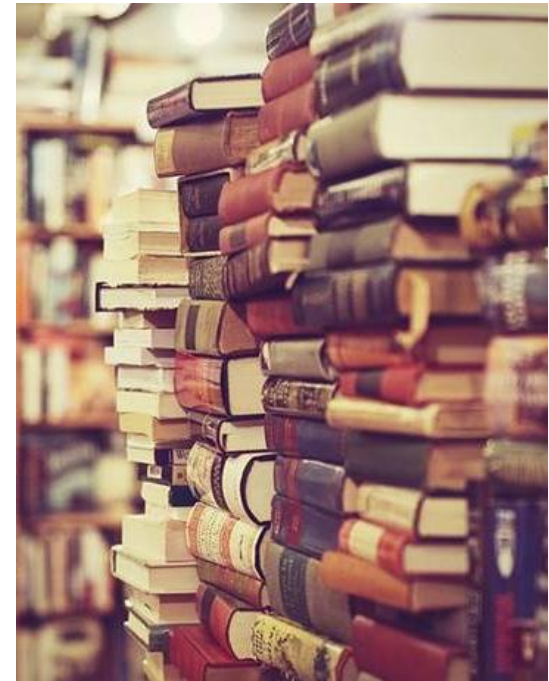
- ◆ $y_* = \arg \max_y p(y | \mathbf{x}, \mathbf{y}, \mathbf{X}, \sigma^2, \tau^2)$
- ◆ Inference problem: what is the most likely state of node y , given observed nodes y_1, \dots, y_N ?

Agenda

- Graphical models: syntax and semantics.
- Inference in graphical models (exact, approximate)
- **Graphical models in machine learning**
 - ◆ Recap: Bayesian linear regression
 - ◆ **Latent Dirichlet Allocation**
 - ◆ Hidden Markov models

Topic Models

- „Topic models“: class of models for the (unsupervised) analysis of text corpora.
- Given a collection of text documents:
 - ◆ Discover the hidden themes (“topics”) that pervade the collection.
 - ◆ Describe topics by the words that most frequently appear.
 - ◆ Describe documents by the topics they are about.
 - ◆ Inferred annotations can be used to organize, summarize, and make predictions.
- Analysis is unsupervised, exploratory in nature.



Latent Dirichlet Allocation

- We will discuss *latent Dirichlet allocation* (LDA)
 - ◆ Well-funded probabilistic model.
 - ◆ Many practical applications.
 - ◆ Easily expressed as a graphical model.

Primer: Categorical and Dirichlet Distributions

- Let $X \in \{v_1, \dots, v_K\}$ denote a discrete random variable that takes on one of K values.
- The *categorical distribution* over X is given by

$$p(X = v_k) = \theta_k$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^T \in R^K$ is a vector of probabilities, that is, $\sum_{k=1}^K \theta_k = 1$.

- The categorical distribution generalizes the Bernoulli distribution.
- We also write $X \sim \text{Cat}(X | \boldsymbol{\theta})$
- Example: rolling a fair dice, $\boldsymbol{\theta} = (1/6, \dots, 1/6)^T$.

Primer: Categorical and Dirichlet Distributions

- The Dirichlet distribution, given by

$$p(\boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha}) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \theta^{\alpha_k - 1} \quad \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$$

generalizes the Beta-distribution (identical to Beta for $K=2$).

- Reminder: $\text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1 - 1} (1 - \theta)^{\alpha_2 - 1}$
- Conjugate prior: Dirichlet + categorical behaves like Beta + Bernoulli

Primer: Categorical and Dirichlet Distributions

- Recap: Beta is conjugate to Bernoulli

- ◆ If $\theta \sim \text{Beta}(\theta | \alpha_1, \alpha_2)$

$$X_i \sim \text{Bern}(X | \theta)$$

- ◆ Then $p(\theta | X_1, \dots, X_N) = \text{Beta}(\theta | \alpha_1 + n_1, \alpha_2 + n_2)$

Number of $X_i = 0$

Number of $X_i = 1$

- Dirichlet is conjugate to categorical

- ◆ If $\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha})$

$$X_i \sim \text{Cat}(X | \boldsymbol{\theta})$$

- ◆ Then $p(\boldsymbol{\theta} | X_1, \dots, X_N) = \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\alpha} + \mathbf{n})$

$$\mathbf{n} = (n_1, \dots, n_K)^T$$

$n_k =$ number of $X_i = v_k$

Primer: Symmetric Dirichlet Distribution

- Dirichlet distribution „smoothes“ probability estimates towards the prior information in the vector $\alpha \in R^K$.
- Example ($K=3$):

Prior: $\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\theta} \mid (5, 5, 5))$

Observations: $\mathbf{n} = (3, 0, 2)$

MAP estimate: $\boldsymbol{\theta}^* = (0.41, 0.23, 0.36)$

Maximum likelihood estimate: $\boldsymbol{\theta}^* = (0.6, 0, 0.4)$

- We often study *symmetric* Dirichlet distributions parameterized by a single $\alpha \in R$.

$$\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\theta} \mid \alpha) = \text{Dir}(\boldsymbol{\theta} \mid (\alpha, \dots, \alpha))$$

$$\alpha_1 = \alpha_2 = \dots = \alpha_K = \alpha$$

Topic Models: Motivating Example

- Example: Analyze the following five sentences (=documents):
 - (1) I like to eat broccoli and bananas.
 - (2) I ate a banana and spinach smoothie for breakfast.
 - (3) Chinchillas and kittens are cute.
 - (4) My sister adopted a kitten yesterday.
 - (5) Look at this cute hamster munching on a piece of broccoli.
- LDA: automatically discover topics that sentences contain.
- Possible answer:
 - ◆ Sentences (1) and (2): 100% Topic 1.
 - ◆ Sentences (3) and (4): 100% Topic 2.
 - ◆ Sentence 5: 60% Topic 1, 40% Topic 2.
 - ◆ Topic 1: 30% broccoli, 15% bananas, 10% breakfast, 10% munching, ...
 - ◆ Topic 2: 20% chinchillas, 20% kittens, 20% cute, 15% hamster, ...

Topic Models: Motivating Example

- Example: Analyze the following five sentences (=documents):

- (1) I like to eat broccoli and bananas.
- (2) I ate a banana and spinach smoothie for breakfast.
- (3) Chinchillas and kittens are cute.
- (4) My sister adopted a kitten yesterday.
- (5) Look at this cute hamster munching on a piece of broccoli.

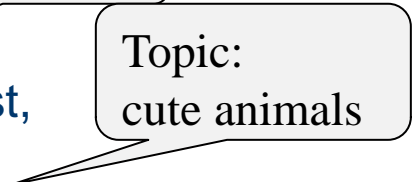
- LDA: automatically discover topics that sentences contain.

- Possible answer:

- ◆ Sentences (1) and (2): 100% Topic 1.
- ◆ Sentences (3) and (4): 100% Topic 2.
- ◆ Sentence 5: 60% Topic 1, 40% Topic 2.
- ◆ Topic 1: 30% broccoli, 15% bananas, 10% breakfast, 10% munching, ...
- ◆ Topic 2: 20% chinchillas, 20% kittens, 20% cute, 15% hamster, ...



Topic: food



Topic:
cute animals

Latent Dirichlet Allocation

- Formalization: documents, words, and topics.
 - ◆ There are D documents, indexed by $d = 1, \dots, D$.
 - ◆ There is a vocabulary of V different words.
 - ◆ Each document contains (up to) N words, denoted by $w_{d,1}, \dots, w_{d,N}$.
 - ◆ There are K topics, indexed by $k = 1, \dots, K$.

- Each topic k is described by a categorical distribution over words, represented by a parameter vector $\beta_k \in R^V$.

Vocabulary = {bananas, broccoli, breakfast, chincillas, cute, hamster, ... }

$$\beta_1 = (0.15 \quad 0.3 \quad 0.1 \quad 0.01 \quad 0.01 \quad 0.01 \quad \dots)^T \in R^V$$

$$\beta_2 = (0.01 \quad 0.01 \quad 0.01 \quad 0.2 \quad 0.2 \quad 0.15 \quad \dots)^T \in R^V$$

Latent Dirichlet Allocation

- Formalization: documents, words, and topics.
 - ◆ There are D documents, indexed by $d = 1, \dots, D$.
 - ◆ There is a vocabulary of V different words.
 - ◆ Each document contains (up to) N words, denoted by $w_{d,1}, \dots, w_{d,N}$.
 - ◆ There are K topics, indexed by $k = 1, \dots, K$.
- Each document d is described by a categorical distribution over topics, represented by a parameter vector $\theta_d \in R^K$.

$$\theta_1 = (1.0 \quad 0.0) \in R^K$$

$$\theta_2 = (1.0 \quad 0.0) \in R^K$$

$$\theta_3 = (0.0 \quad 1.0) \in R^K$$

$$\theta_4 = (0.0 \quad 1.0) \in R^K$$

$$\theta_5 = (0.6 \quad 0.4) \in R^K$$

LDA: Generative Process

- Latent Dirichlet Allocation is a generative model, defining a generative process for the words appearing in the D documents.

- For topics $k = 1, \dots, K$:

Dirichlet prior

- ◆ Draw the categorical distribution over words, $\beta_k \sim \text{Dir}(\beta|\eta)$.

- For documents $d = 1, \dots, D$:

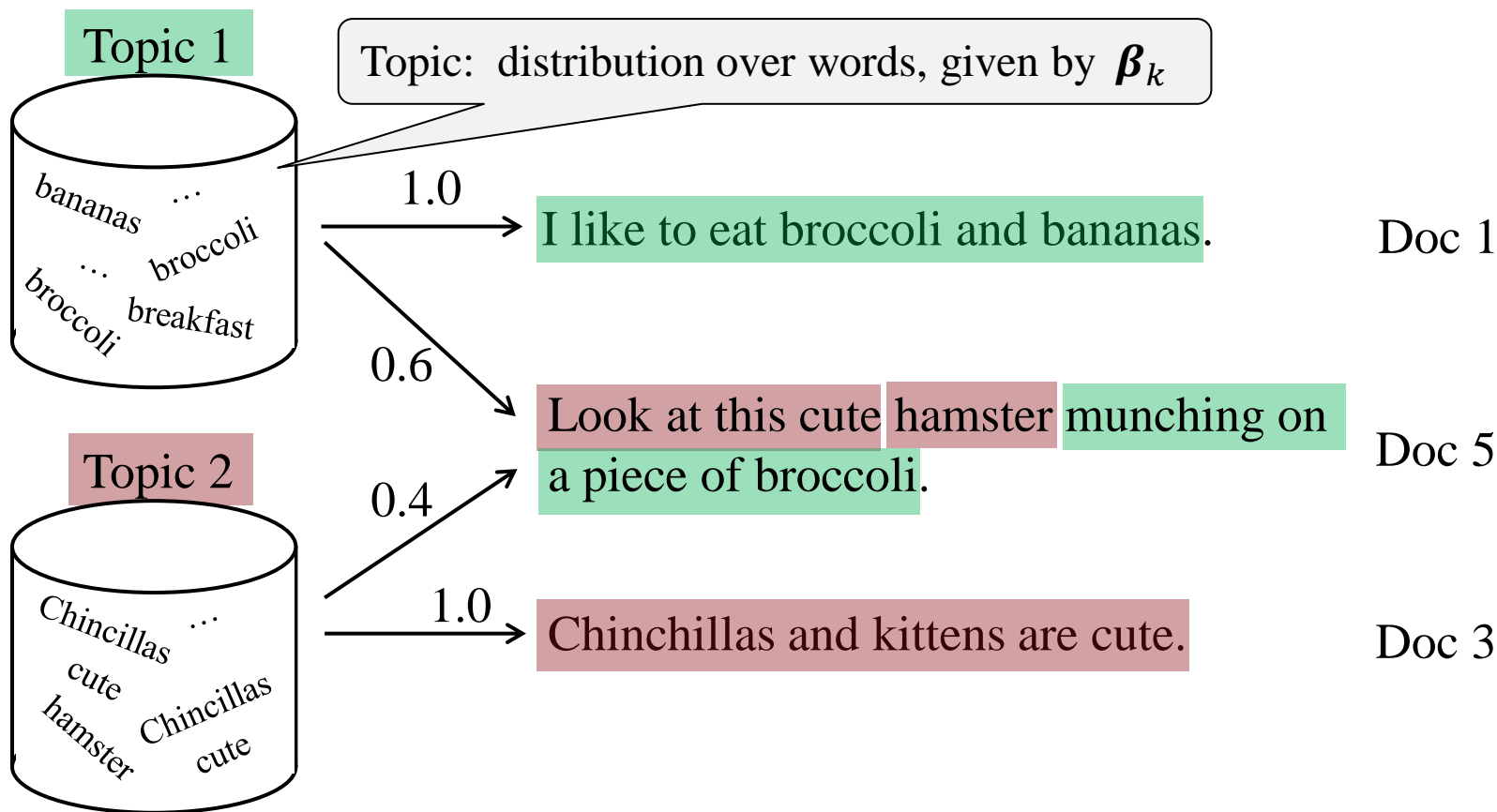
Dirichlet prior

- ◆ Draw the categorical distribution over topics, $\theta_d \sim \text{Dir}(\theta|\alpha)$.
- ◆ For each word $w_{d,1}, \dots, w_{d,N}$:
 - ★ Draw a topic $Z_{d,n} \sim \text{Cat}(Z|\theta_d)$ for the word at position n .
 - ★ Draw the word $w_{d,n} \sim \text{Cat}(w|\beta_{Z_{d,n}})$

Distribution over words
in the selected topic.

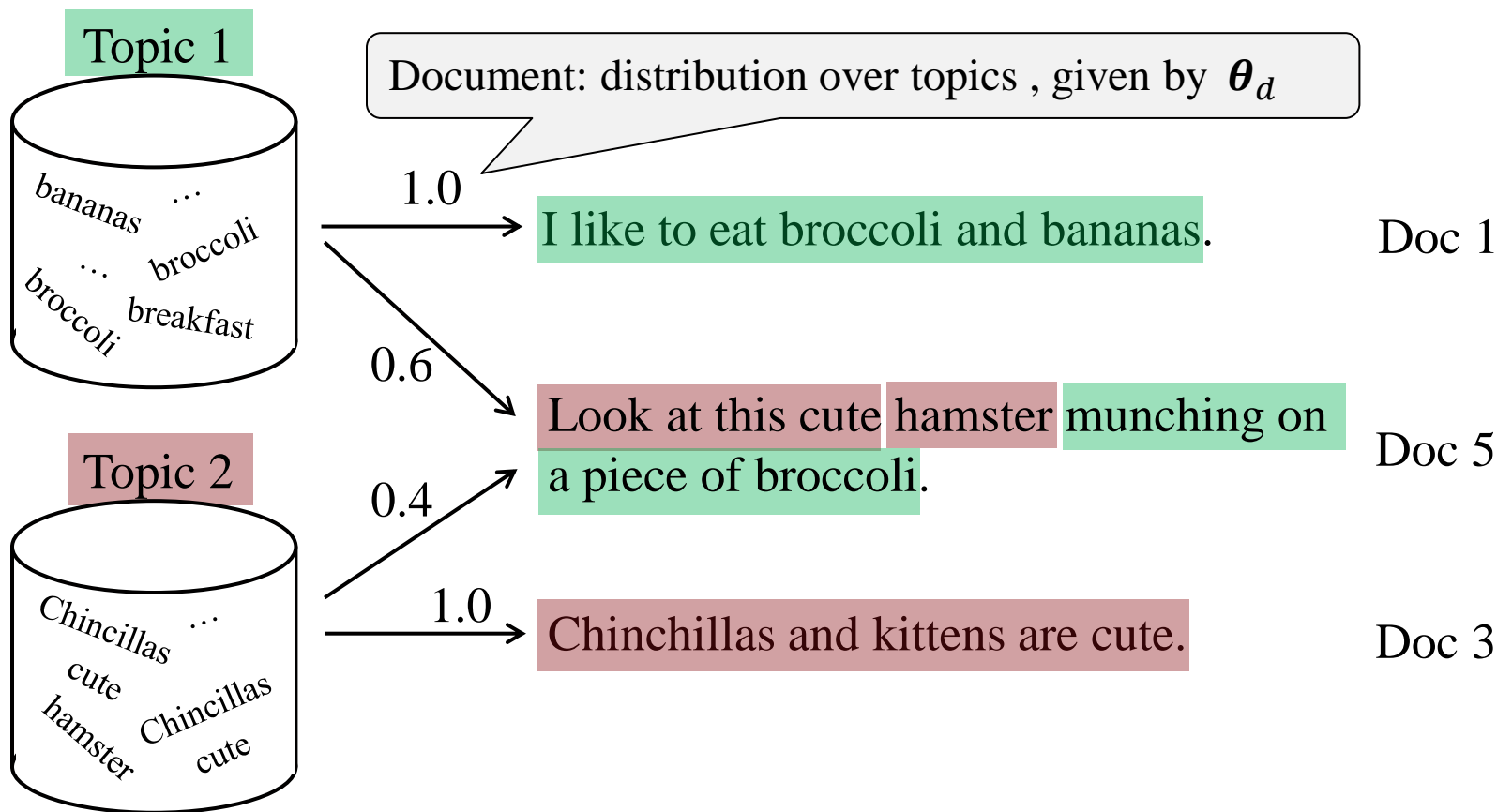
LDA: Generative Process

- Visualization of generative process for documents.



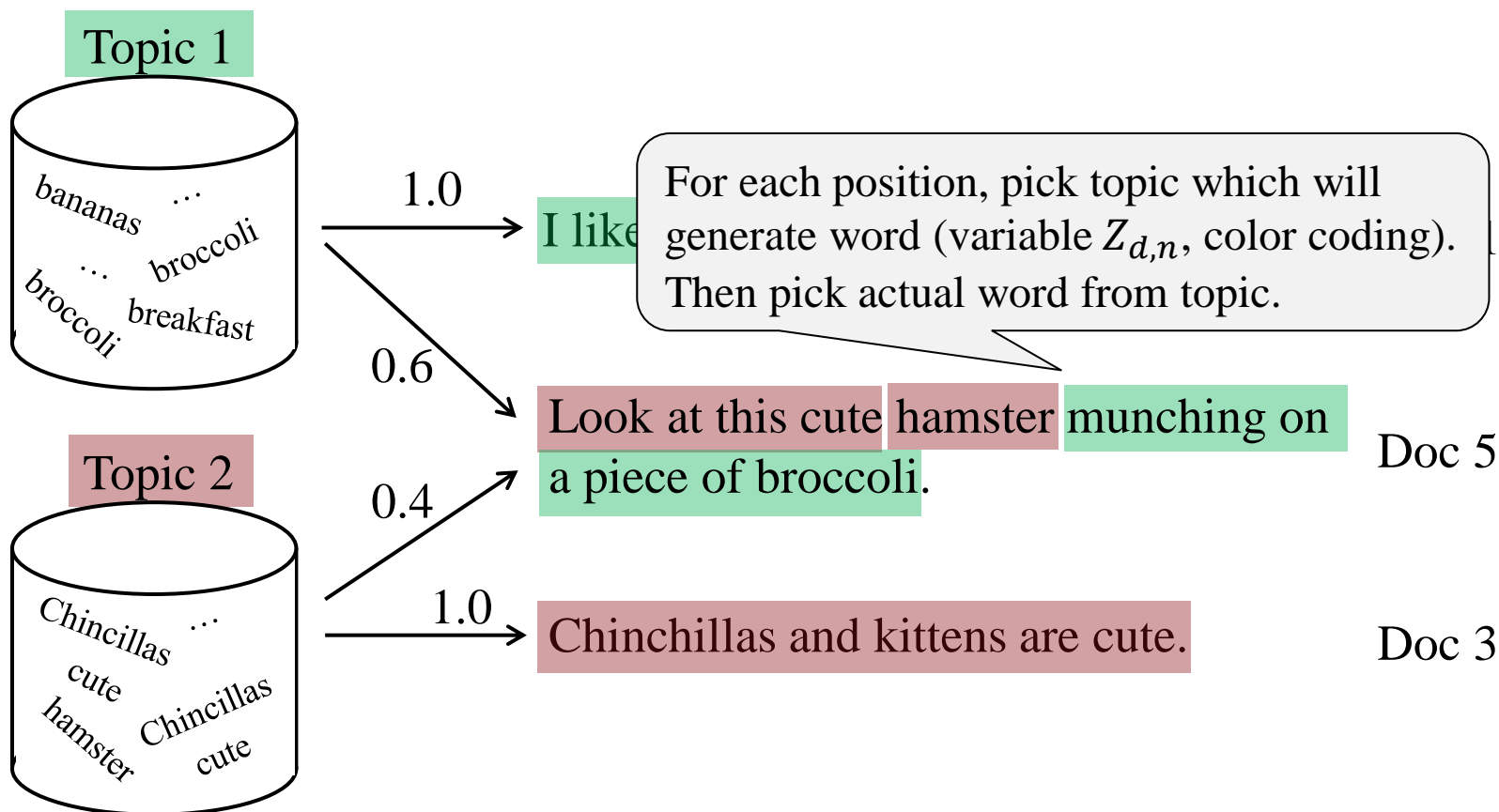
LDA: Generative Process

- Visualization of generative process for documents.



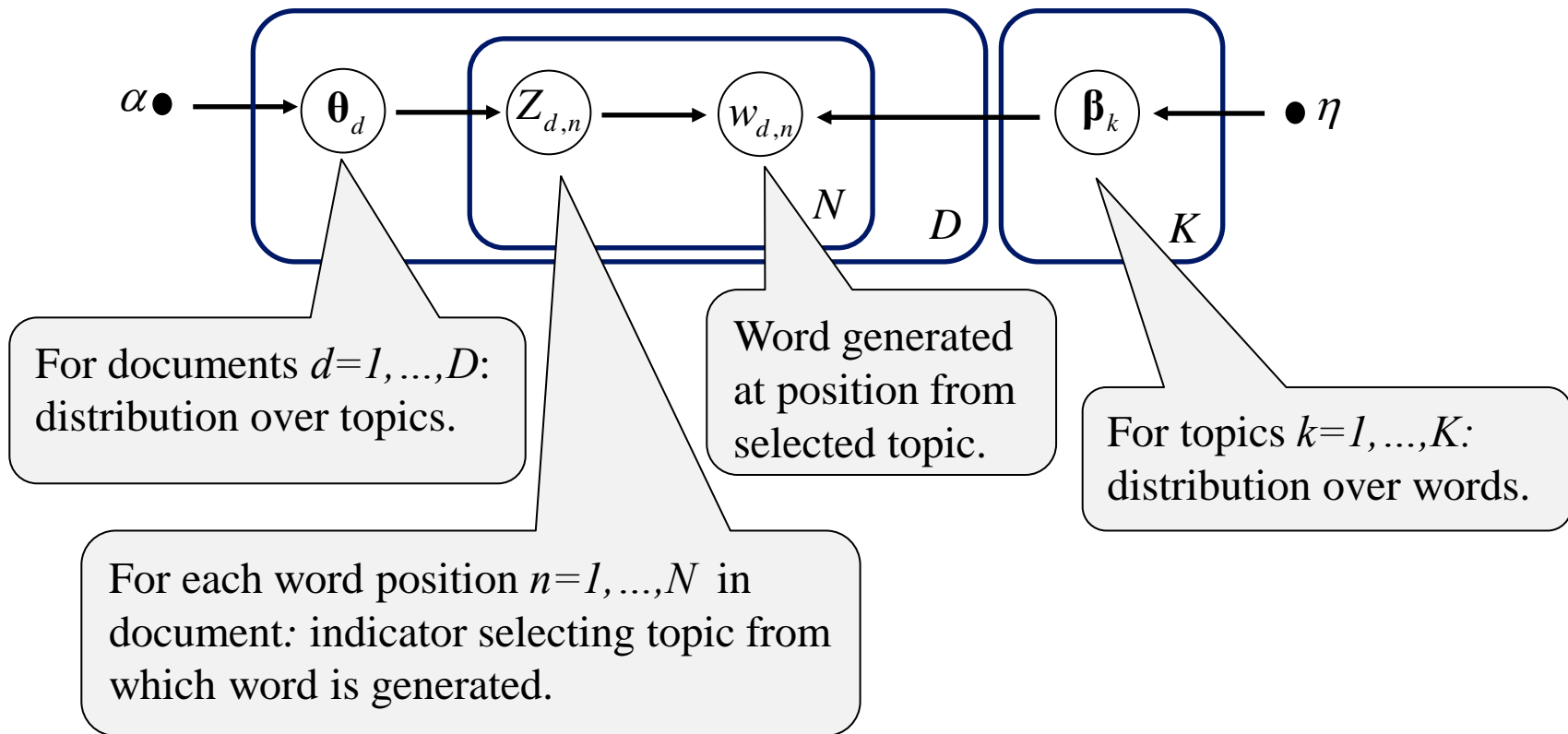
LDA: Generative Process

- Visualization of generative process for documents.



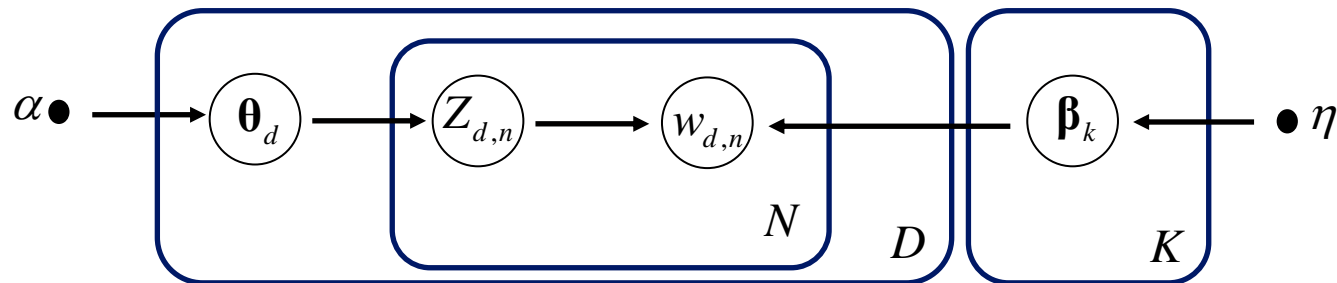
LDA: Graphical Model

- LDA as a graphical model (nested plate notation).

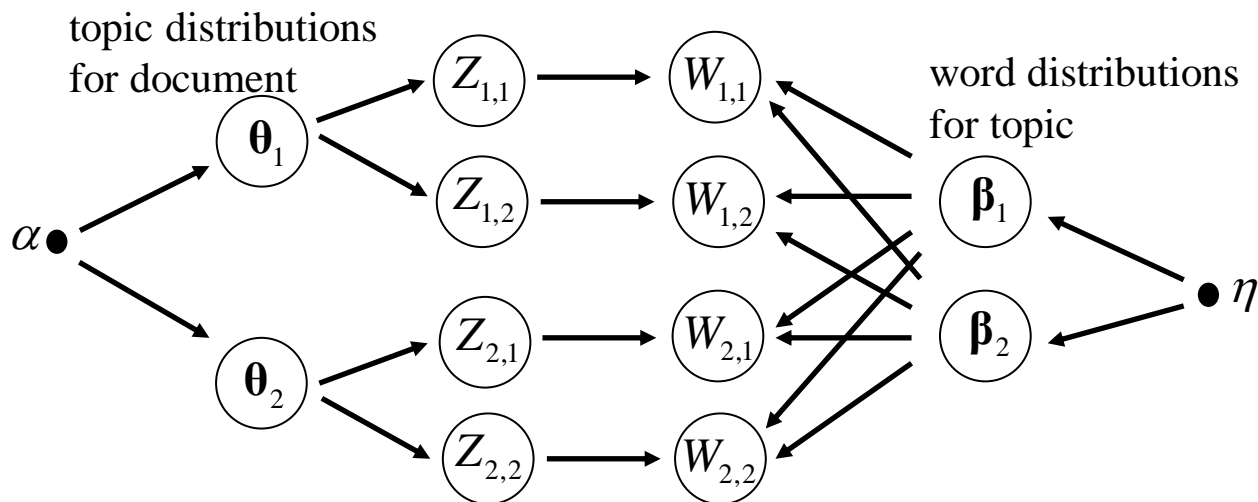


LDA: Graphical Model

- LDA as a graphical model (nested plate notation).

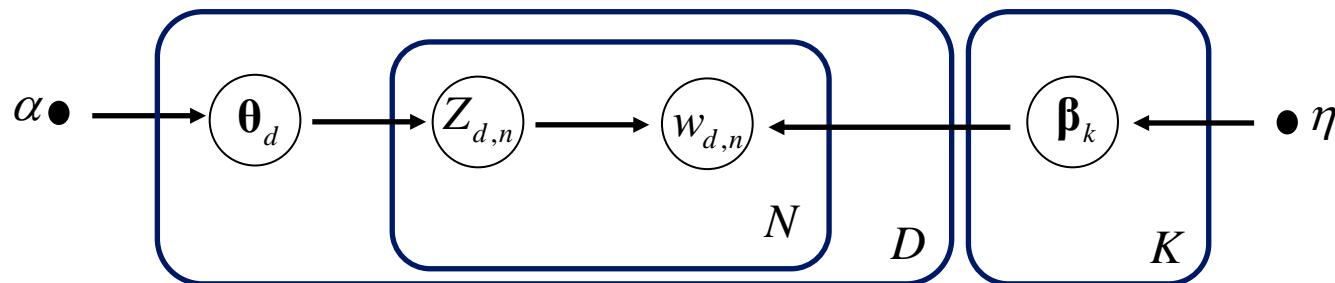


- Unrolled graphical model for small example ($D=N=K=2$):



LDA: Graphical Model

- LDA as a graphical model (nested plate notation).



- Collect all θ_d in stacked vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_D)^T$.
- Collect all β_k in stacked vector $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^T$.
- Collect all $Z_{d,n}$ in matrix $\mathbf{Z} \in R^{D \times N}$.
- Collect all $w_{d,n}$ in matrix $\mathbf{W} \in R^{D \times N}$.
- Joint distribution:

$$p(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{Z}, \mathbf{W} | \eta, \alpha) = \left(\prod_{k=1}^K p(\beta_k | \eta) \right) \left(\prod_{d=1}^D p(\theta_d | \alpha) \prod_{n=1}^N p(Z_{d,n} | \theta_d) p(w_{d,n} | \boldsymbol{\beta}, Z_{d,n}) \right)$$

LDA For Text Analysis

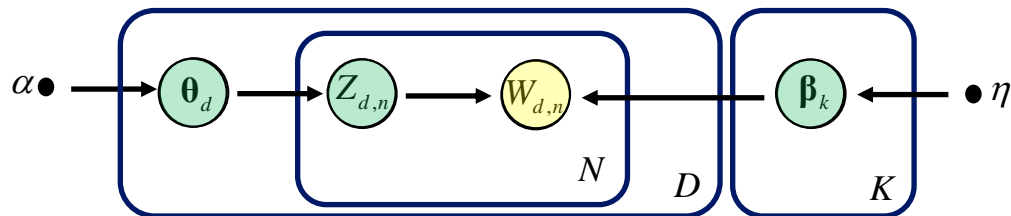
- Although LDA is a generative model, it is not actually good at generating natural language texts.

- LDA is a bag of words model:
 - ◆ To generate a word in a document, pick a random topic ($Z_{d,n}$ variable) and generate a word from that topic ($w_{d,n}$).
 - ◆ No sequential information: likely to generate texts like „breakfast broccoli bananas munching“ or „kitten cute hamster look“.

- **Practical application is in text analysis:**
 - ◆ Document collection is given. Discover hidden topics present in given document collection.
 - ◆ Annotate documents with topics.
 - ◆ This actually works well.

LDA For Text Analysis: Inference Problem

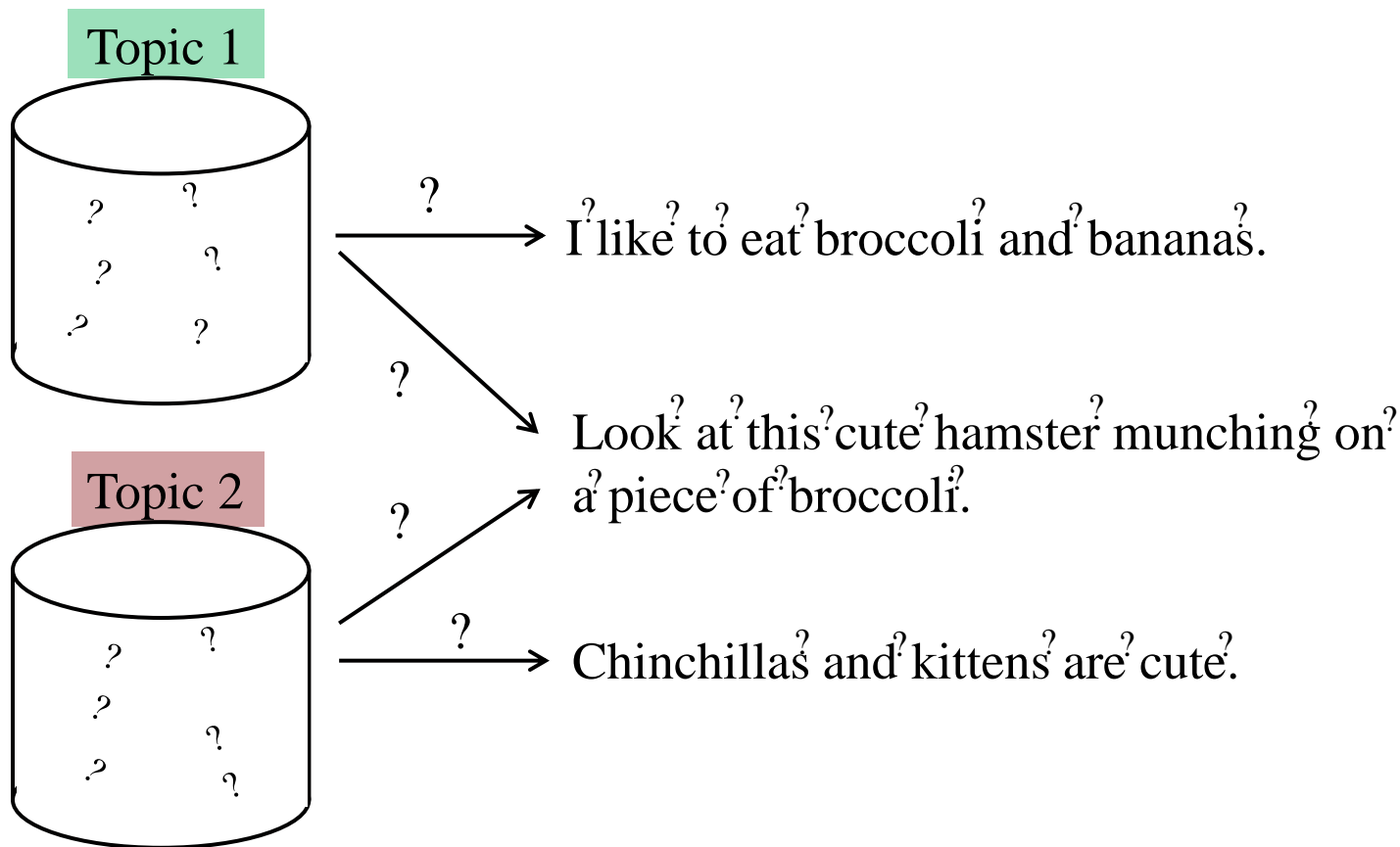
- Problem setting: given document collection, infer topic structure.
 - ◆ **Given:** collection of D documents with N words each, represented by variables $w_{d,n}$ ($d = 1, \dots, D, n = 1, \dots, N$).
 - ◆ **Infer:**
 - ★ topic distribution for each document (variables θ_d),
 - ★ word distribution for each topic (variables β_k)
 - ★ which topic has generated a word ($Z_{d,n}$).
- Inference problem: compute $p(\theta, \beta, Z | W)$.



- All variables except the word information W are *latent variables*, that is, they are never observed.

LDA For Text Analysis: Inference Problem

- Need to infer all other variables given word information.



LDA: Inference Algorithms

- Computing $p(\boldsymbol{\theta}, \boldsymbol{\beta}, \mathbf{Z}|\mathbf{W})$ means solving an inference problem in the graphical model representing LDA.

- Different approximate inference algorithms have been studied:
 - ◆ Gibbs sampler (Pritchard et al., 2000)
 - ◆ Collapsed Gibbs sampler (Griffiths & Steyvers, 2004)
 - ◆ Variational methods (Blei et al., 2003)
 - ◆ Expectation propagation (Minka and Lafferty, 2002)
 - ◆ ...

- Sampling-based approaches very popular.

LDA: Gibbs Sampling

- Recap: Gibbs sampling resamples each variable given all other variables.
- Sample discrete variables $Z_{d,n}$, $w_{d,n}$: as discussed.
- Sample continuous variables θ_d (topic distribution for document):

D-separation

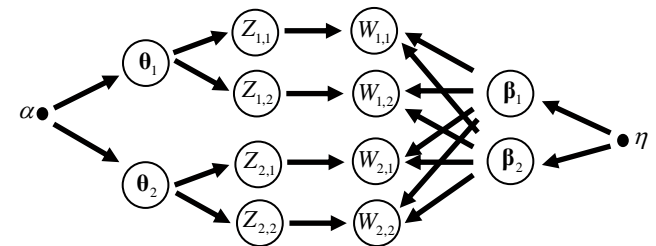
$$p(\theta_1 | \theta_2, \beta, \mathbf{W}, \mathbf{Z}, \eta, \alpha) = p(\theta_1 | Z_{1,1}, Z_{1,2}, \alpha)$$

$$= \text{Dir}(\theta_1 | (n_1 + \alpha, \dots, n_K + \alpha))$$

Conjugate prior:
Posterior Dirichlet.

$$\text{with } n_k = \sum_{n=1}^N \mathbb{I}[Z_{1,n} = k]$$

Count how often Topic k was chosen in Document 1.



LDA: Gibbs Sampling

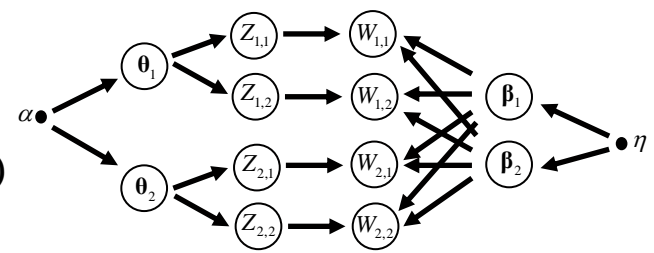
- Sample continuous variables β_k (word distribution for topic):

$$p(\beta_1 | \theta, \beta_2, \mathbf{W}, \mathbf{Z}, \eta, \alpha) = p(\beta_1 | \mathbf{W}, \mathbf{Z}, \eta)$$

D-separation

$$= \text{Dir}(\beta_1 | (m_1 + \eta, \dots, m_v + \eta))$$

Conjugate prior:
Posterior Dirichlet.



$$\text{with } m_v = \sum_{d=1}^D \sum_{n=1}^N \mathbb{I}[w_{d,n} = v \wedge Z_{d,n} = 1]$$

Count how often word v was generated from Topic 1 (in any document).

- We can easily sample θ_d and β_k from their respective Dirichlet posterior.

LDA: Sampling Inference

- We obtain samples $(\boldsymbol{\theta}^t, \boldsymbol{\beta}^t, \mathbf{Z}^t)$ for $t=1, \dots, T$ from Gibbs sampler.
- Values for variables of interest are usually derived by averaging over samples (= expected value under distribution):

$$\boldsymbol{\theta} = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\theta}^t \quad \boldsymbol{\beta} = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\beta}^t \quad \mathbf{Z} = \frac{1}{T} \sum_{t=1}^T \mathbf{Z}^t$$

- In practice, full Gibbs sampling rarely used.
- The popular „collapsed“ Gibbs sampler exploits conjugacy of the priors over $\boldsymbol{\theta}$ and $\boldsymbol{\beta}$ to integrate out these variables and only run the Gibbs sampler over the topic assignments \mathbf{Z} .

LDA: Example

- 100-Topic model trained on a corpus of 16,000 *Associated Press* documents.

| “Arts” | “Budgets” | “Children” | “Education” |
|---------|-----------|------------|-------------|
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |

Top words
for 4 topics.

Document
about these
topics

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

LDA: Summary and Remarks

- LDA: probabilistic model for discovering hidden topics in a document collection, and describing documents by these topics.
- Generative model:
 - ◆ A topic is a distribution over words.
 - ◆ A document is a distribution over topics.
 - ◆ For each position in document, pick topic and generate word.
- All variables except words are latent, inferred during inference.
- For simplicity, we have defined the model assuming that the number N of words is the same in each document
 - ◆ Straightforward to generalize to N differing over documents.
 - ◆ Generate N from model (Poisson distribution).

Agenda

- Graphical models: syntax and semantics.
- Inference in graphical models (exact, approximate)
- **Graphical models in machine learning**
 - ◆ Recap: Bayesian linear regression
 - ◆ Latent Dirichlet allocation
 - ◆ **Hidden Markov models**

Hidden Markov Model: Probabilistic Automaton with Hidden States

- Hidden Markov model: probabilistic model for sequences.
- Temporal view: discrete time steps $1, \dots, T$ (= sequence elements).
- Models a probabilistic automaton that takes on a state from a finite set of states at each point in time.
- At each point in time, the automaton probabilistically changes into a novel state, based on the current state.

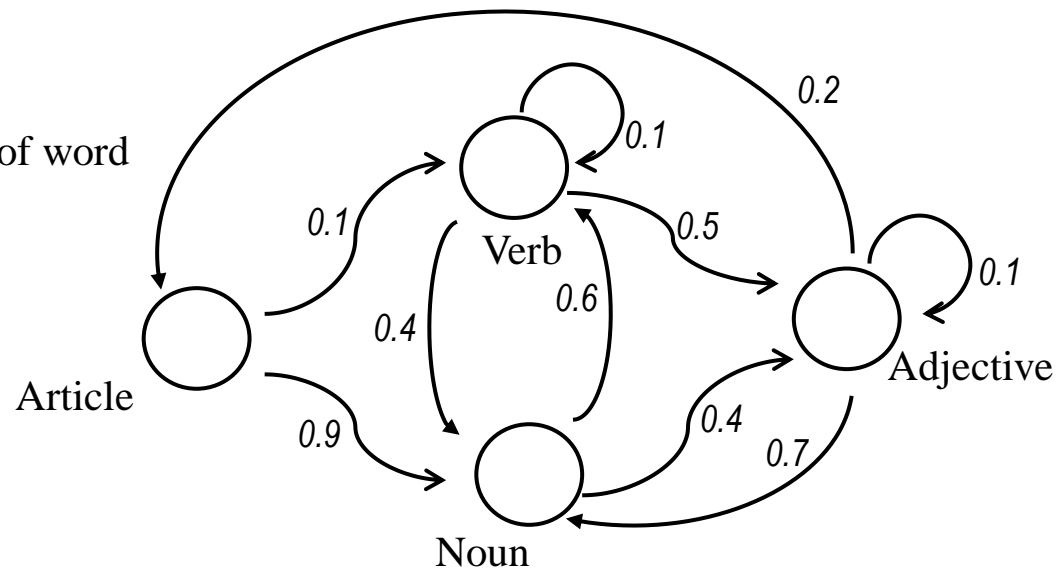
Example:

Model for natural language

Time steps = words in text

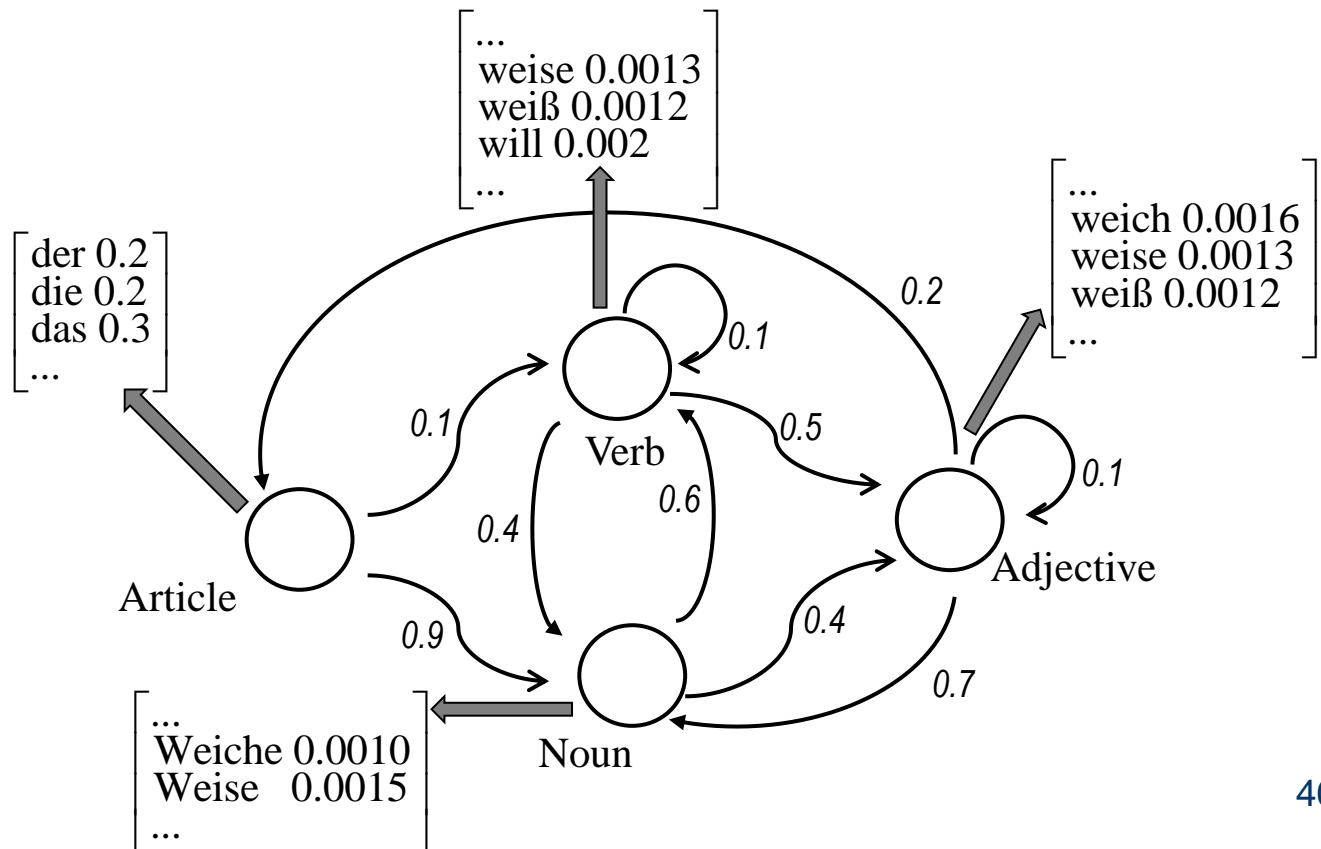
States = grammatical category of word

Probabilities for
state transitions



Hidden Markov Model: Probabilistic Automaton with Hidden States

- Sequence of states is not directly observable.
- Instead: states emit observations, these form the observable sequence.
- Distribution over possible observations depends on the current state.



Formalization: States

- State at time t : random variable $q_t \in \{1, \dots, N\}$
- Automaton model:
 - ◆ Distribution over random initial state: $p(q_1)$
 - ◆ Distribution over next state given current state: $p(q_t | q_{t-1})$
- This results in joint distribution over all states:

$$p(q_1, \dots, q_T) = p(q_1) \prod_{t=2}^T p(q_t | q_{t-1})$$

- As a graphical model:



- „Markov“-Assumption: $p(q_t | q_1, \dots, q_{t-1}) = p(q_t | q_{t-1})$

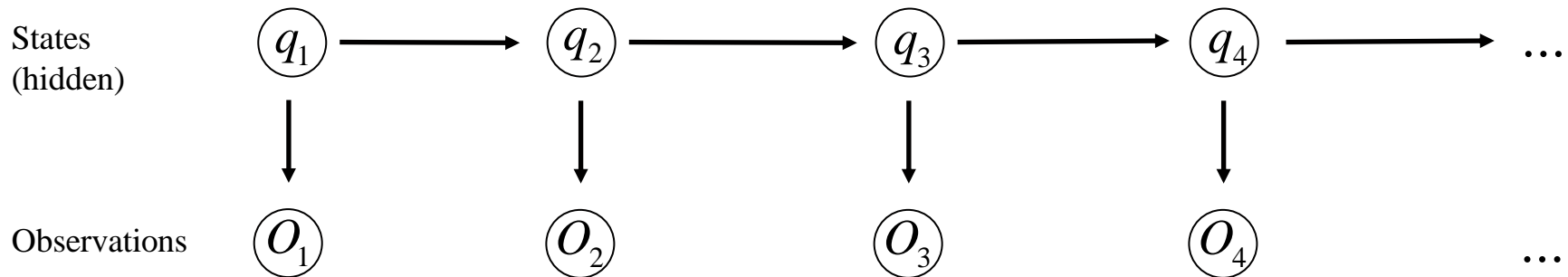
Formalization: Observations

- Observation at time t : random variable $O_t \in \{o_1, \dots, o_M\}$.
- Automaton model: observation is generated depending on current state, from distribution $p(O_t | q_t)$.

- Joint distribution over all random variables:

$$p(q_1, \dots, q_T, O_1, \dots, O_T) = p(q_1) p(O_1 | q_1) \prod_{t=2}^T p(q_t | q_{t-1}) p(O_t | q_t)$$

- As graphical model („Hidden Markov Model“ or „HMM“)



HMM Parameterization

- To define a hidden Markov model, we need to specify the following distributions:

$p(q_1)$ Distribution over initial states

$p(q_t | q_{t-1})$ Distribution over state transitions (independent of t)

$p(O_t | q_t)$ Distribution over observations (independent of t)

- Notation for the corresponding probability values:

Initial state probabilities: $p(q_1 = j) = \pi_j$ Vector $\pi \in \mathbb{R}^N$

Transition probabilities: $p(q_t = j | q_{t-1} = i) = a_{ij}$ Matrix $A \in \mathbb{R}^{N \times N}$

Observation probabilities: $p(O_t = o_m | q_t = i) = b_i(o_m)$ Matrix $B \in \mathbb{R}^{M \times N}$

- A hidden Markov mode is defined by the triple $\lambda = (A, B, \pi)$.

HMM Problem Settings

- Three basic problem settings for hidden Markov models:
- 1. Likelihood of an observation sequence: How likely is a sequence of observations O_1, O_2, \dots, O_T given a model λ ?

$$\text{Compute } p(O_1, \dots, O_T | \lambda)$$

- 2. Most probable state, given observations:
 - a) Compute $\arg \max_{q_t} p(q_t | O_1, \dots, O_T, \lambda)$
 - b) Compute $\arg \max_{q_1, \dots, q_T} p(q_1, \dots, q_T | O_1, \dots, O_T, \lambda)$
- 3. Given several observation sequences, find a model that best explains the data:

$$\text{Compute } \arg \max_{\lambda} p(\{(O_1, \dots, O_T), \dots\} | \lambda)$$

1. Likelihood of an Observation Sequence

- Likelihood of an observation sequence: How likely is an observation sequence O_1, O_2, \dots, O_T given a model λ ?
- Sum rule:

$$p(O_1, \dots, O_T | \lambda) = \sum_{q_1} \dots \sum_{q_T} p(O_1, \dots, O_T, q_1, \dots, q_T | \lambda)$$

Exponential
time

- Goal: polynomial-time algorithm.
- Solution:
 - ◆ Forward-Backward algorithm: dynamic programming.
 - ◆ Forward-Backward is special case of (general) message passing.
 - ◆ Also solves Problem 2.a)

Forward-Backward Algorithm

- Define auxiliary variables

$$\alpha_t(i) = p(O_1, \dots, O_t, q_t = i \mid \lambda)$$

$N \cdot T$ variables overall

$$\beta_t(i) = p(O_{t+1}, \dots, O_T \mid q_t = i, \lambda)$$

$N \cdot T$ variables overall

- Theorem:

$$\alpha_1(i) = \pi_i b_i(O_1)$$

$$\alpha_{t+1}(i) = \left(\sum_{j=1}^N \alpha_t(j) a_{ji} \right) b_i(O_{t+1})$$

Allows recursive computation of all $\alpha_t(i)$
for $t = 1, \dots, T$ and $j = 1, \dots, N$ in time $O(N^2T)$

$$\beta_T(i) = 1$$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)$$

Allows recursive computation of all $\beta_t(i)$
for $t = T, \dots, 1$ and $i = 1, \dots, N$ in time $O(N^2T)$

Forward-Backward Algorithm

- If we have computed all $\alpha_t(i)$ and $\beta_t(i)$:

Solution problem 1:

$$\begin{aligned} p(O_1, \dots, O_T | \lambda) &= \sum_{i=1}^N p(O_1, \dots, O_T, q_T = i | \lambda) && \text{sum rule} \\ &= \sum_{i=1}^N \alpha_T(i) \end{aligned}$$

Total time: $O(N^2T)$

Forward-Backward Algorithm

- If we have computed all $\alpha_t(i)$ and $\beta_t(i)$:

Solution Problem 2.a):

$$\begin{aligned} p(q_t = i | O_1, \dots, O_T, \lambda) &= \dots \\ &= \frac{\alpha_t(i)\beta_t(i)}{p(O_1, \dots, O_T | \lambda)} \end{aligned}$$

Total time: $O(N^2T)$

Outlook

- Problem 2.b): Solution with Viterbi algorithm (same idea as in Forward-Backward)
- Problem 3: Solution with Baum-Welch algorithm
 - ◆ Instance of EM-Algorithm (see also Gaussian Mixture Models)
 - ◆ Makes use of Forward-Backward in E-step.
- Alternative approaches for modeling sequential data: discriminative models
 - ◆ HM-SVM
 - ◆ Conditional Random Fields
- More details in the lecture „speech technology“.

Applications HMMs

- Part-of-speech tagging in natural language texts
 - ◆ Hidden states correspond to grammatical categories (article, verb, noun, etc).
 - ◆ Observations are words in text.
 - ◆ Goal: assign grammatical categories to words (= find most likely hidden state sequence).
- Speech recognition: acoustic model
 - ◆ Hidden states are spoken words.
 - ◆ Observation is the acoustic signal.
 - ◆ Goal: reconstruct spoken words from acoustic signal.
 - ◆ Partially superseded by deep learning.
- Bioinformatics
 - ◆ Localization or annotation of genes.
 - ◆ Usually extensions of the basic hidden Markov model.