

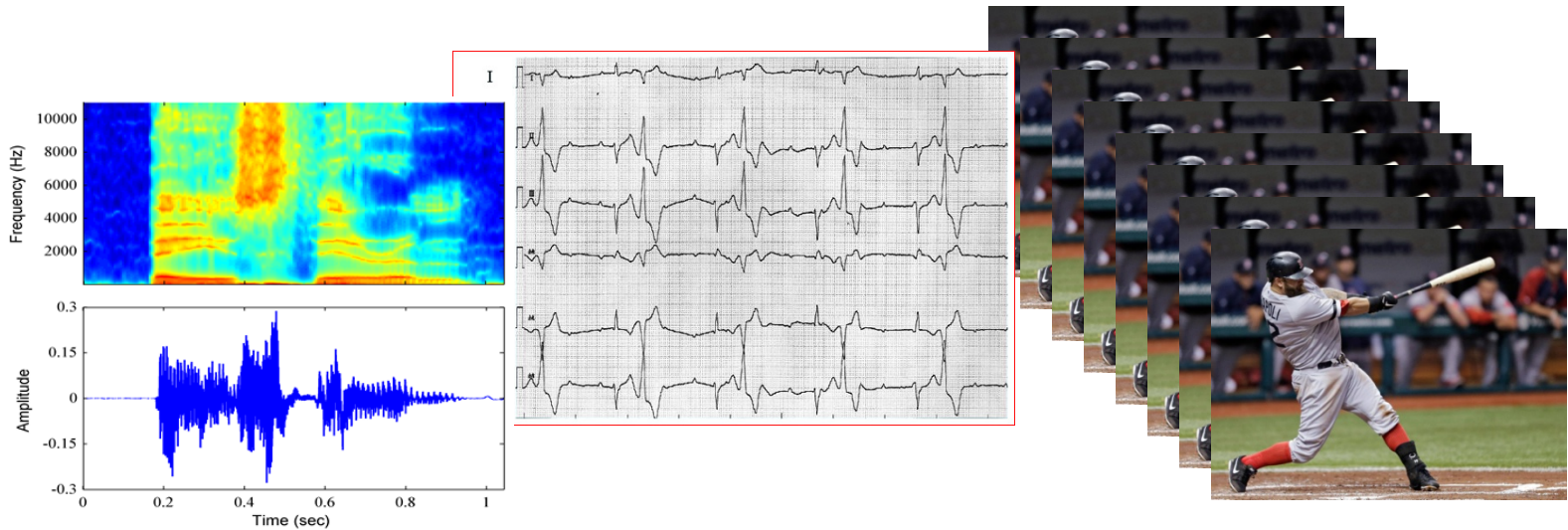
Sequential Data Modeling

Tomoki Toda
Graham Neubig
Sakriani Sakti

Augmented Human Communication Laboratory
Graduate School of Information Science

Course Goals

The aim of this course is to **learn basic knowledge of sequential data modeling techniques** that can be applied to sequential data such as **speech signals**, **biological signals**, **videos of moving objects**, or **natural language text**. In particular, it will focus on deepening knowledge of methods based on **probabilistic models**, such as hidden Markov models or linear dynamical systems.



Credits and Grading

- 1 credit course
- Score will be graded by
 - Assignment report in every class
- Prerequisites
 - Fundamental Mathematics for Optimization (最適化数学基礎)
 - Calculus (微分積分学)
 - Basic Data Analysis (データ解析基礎)

Materials

- Textbook
 - There is no textbook for this course.
- Lecture slides
 - Handout will be distributed in each class.
 - PDF slides are available from
<http://ahclab.naist.jp/lecture/2016/sdm/index.html>
(internal access only)
- Reference materials
 - C.M. Bishop: Pattern Recognition and Machine Learning, Springer Science + Business Media, LLC, 2006
 - C.M. ビショップ (著)、元田、栗田、樋口、松本、村田 (訳): パターン認識と機械学習 上・下、シュプリンガー・ジャパン、2008

Office Hours

- NAIST Lecturers: Graham Neubig, Sakriani Sakti
Augmented Human Communication Laboratory
- Office: B714
- Office hour: by appointment (send an email first)
- Email: neubig@is.naist.jp, ssakti@is.naist.jp

- Other Contact
 - Tomoki Toda
 - Email: tomoki@icts.nagoya-u.ac.jp

TA Members Email: **sdm2016@is.naist.jp**

- Rui Hiraoka hiraoka.rui.hj9@is.naist.jp
- Yoko Ishikawa ishikawa.yoko.io5@is.naist.jp

Hiraoka-kun



Ishikawa-san



Schedule

- 1st slot on every **Friday 9:20-10:50** in room **L1**

June

Sun	Mon	Tue	Wed	Thu	Fri	Sat
			1	2	3	4
5	6	7	8	9	10	11
12	13	14	15	16	17	18
19	20	21	22	23	24	25
26	27	28	29	30		

July/August

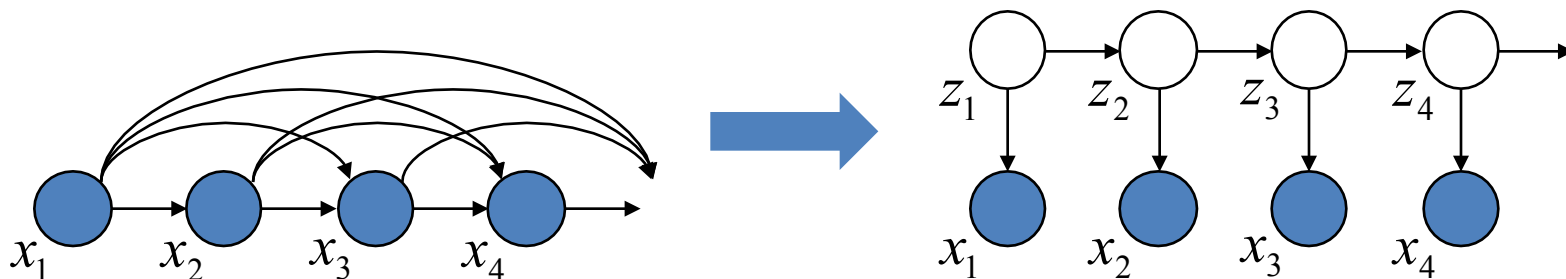
Sun	Mon	Tue	Wed	Thu	Fri	Sat
					1	2
3	4	5	6	7	8	9
10	11	12	13	14	15	16
17	18	19	20	21	22	23
24	25	26	27	28	29	30
31	1					

Syllabus

Date	Course description	Lecturer
6/03	Basics of sequential data modeling 1	Graham Neubig
6/10	Basics of sequential data modeling 2	Graham Neubig
6/17	Discrete latent variable models 1	Tomoki Toda
6/24	Discrete latent variable models 2	Tomoki Toda
7/1	Continuous latent variable models 1	Tomoki Toda
7/15	Discriminative models for sequential labeling 1	Sakriani Sakti
7/29	Continuous latent variable models 2	Tomoki Toda
8/1	Discriminative models for sequential labeling 2	Sakriani Sakti

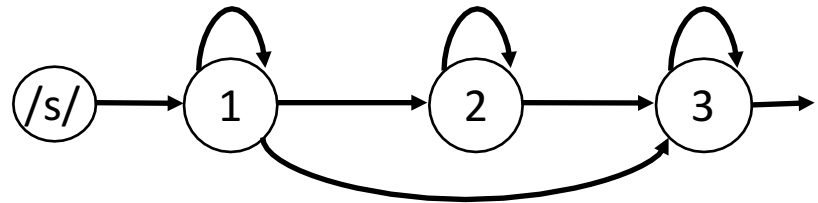
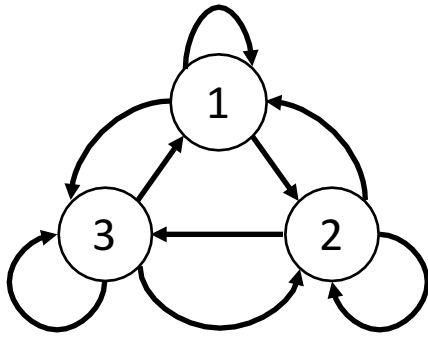
1st and 2nd Classes (6/03 and 6/10)

- Lecturer: Graham Neubig
- Contents: Basics of sequential data modeling
 - Markov process
 - Latent variables
 - Mixture models
 - Expectation-maximization (EM) algorithm



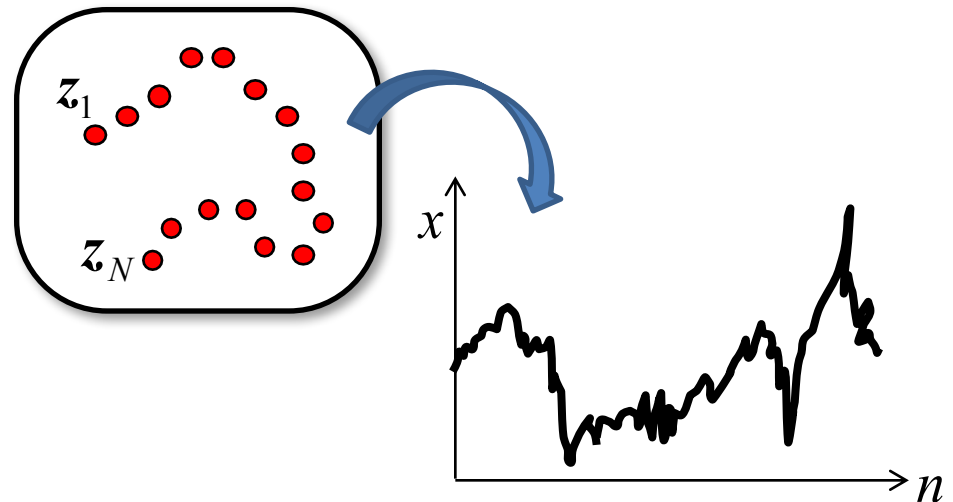
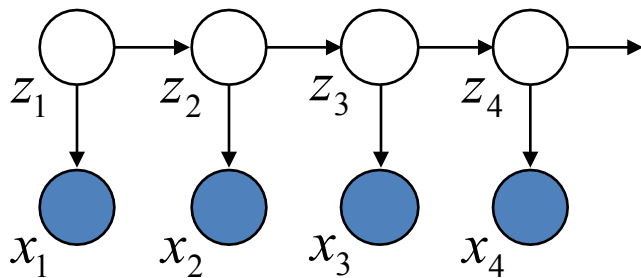
3rd and 4th Classes (6/17 and 6/24)

- Lecturer: Tomoki Toda
- Contents: Discrete latent variable models
 - Hidden Markov models
 - Forward-backward algorithm
 - Viterbi algorithm
 - Training algorithm



5th and 7th Classes (7/1 and 7/29)

- Lecturer: Tomoki Toda
- Contents: Continuous latent variable models
 - Factor analysis
 - Linear dynamical systems
 - Prediction and update
 - (Training algorithm)



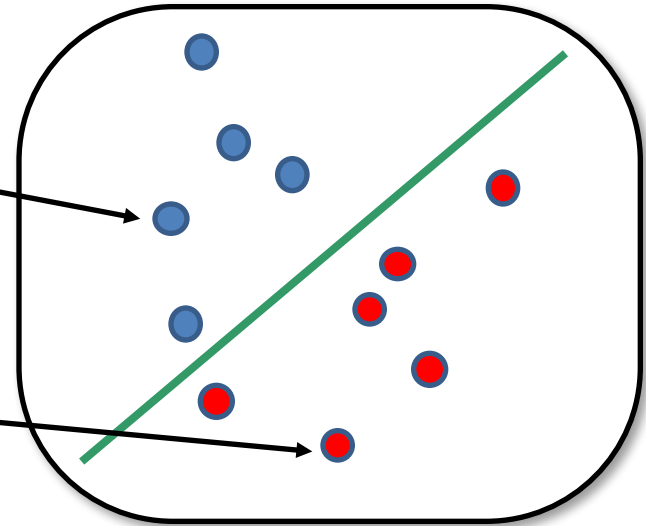
6th and 8th Classes (7/15 and 8/1)

- Lecturer: Sakriani Sakti
- Contents: Discriminative models for sequential labeling
 - Structured perceptron
 - Conditional Random Fields
 - Training algorithm

「今日は晴れた。」

今日／は／晴れ／だ／。

今／日／は／晴れ／だ／。



Sequential Data Modeling

1st class

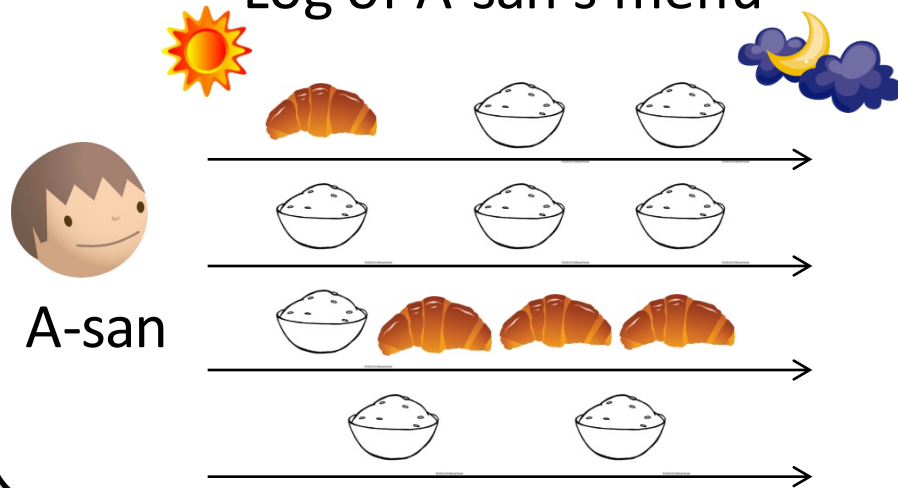
“Basics of sequential data modeling 1”

Graham Neubig

Augmented Human Communication Laboratory
Graduate School of Information Science

Question



Log of A-san's menu



One day, someone ate the following menu.

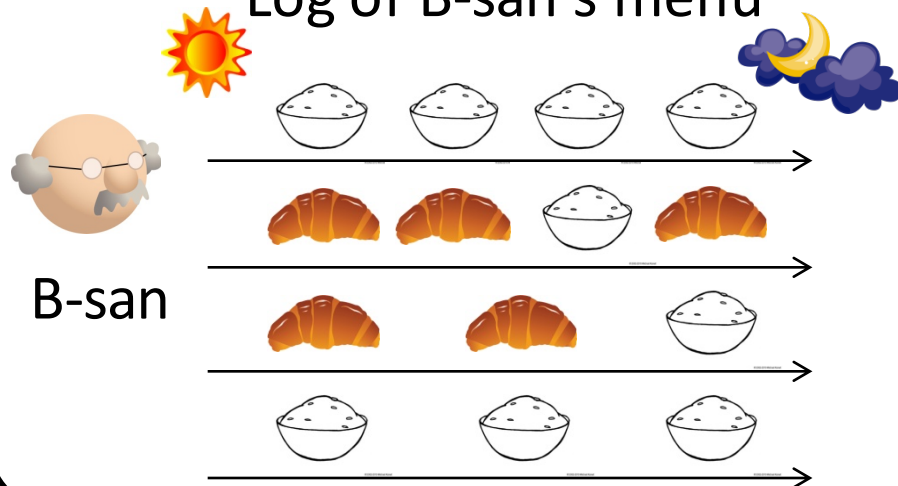


Q1. A-san or B-san?

Q2. If this is A-san's menu, which is "?",  

Q3. ...
⋮

Log of B-san's menu



After this class, you can answer these questions!

Sequential Data

- Data examples
 - Time series (speech, actions, moving objects, exchange rates, ...)
 - Character strings (word sequence, symbol string, ...)
- Various lengths of data
 - E.g.,
 - Data sample 1 (length = 5): { 1, 0, 1, 1, 0 }
 - Data sample 2 (length = 8): { 1, 1, 1, 0, 1, 1, 0, 0 }
 - Data sample 3 (length = 3): { 0, 0, 1 }
 - Data sample 4 (length = 6): { 0, 1, 0, 1, 1, 0 }
- Probabilistic approach to modeling sequential data
 - Consistent framework for the quantification and manipulation of uncertainty
 - Effective for dealing with real data

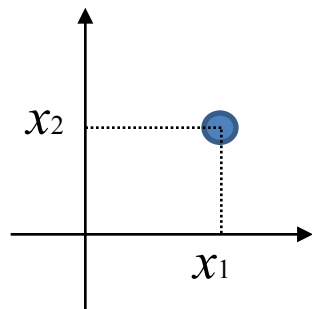
How to Represent Sequential Data?

- A sequential data sample is represented in a high-dimensional space (“# of dimensions” = “length of the sequential data sample”).

Examples of sequential data:

Length = 2

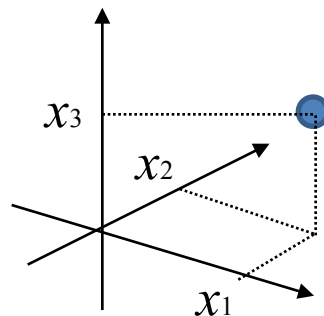
$$\left\{ \begin{array}{cc} x_1 & x_2 \\ n=1 & n=2 \end{array} \right\}$$



Represented by
2-dimensional vector

Length = 3

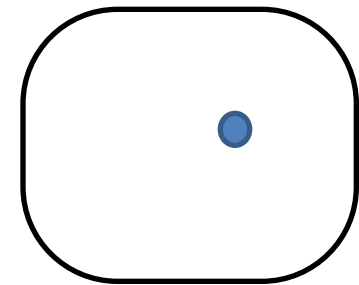
$$\left\{ \begin{array}{ccc} x_1 & x_2 & x_3 \\ n=1 & n=2 & n=3 \end{array} \right\}$$



Represented by
3-dimensional vector

Length = N

$$\left\{ \begin{array}{cccccc} x_1 & x_2 & x_3 & \cdots & x_N \\ n=1 & n=2 & n=3 & \cdots & n=N \end{array} \right\}$$



Represented by
 N -dimensional vector

We need to model probability distribution in these high-dimensional spaces!

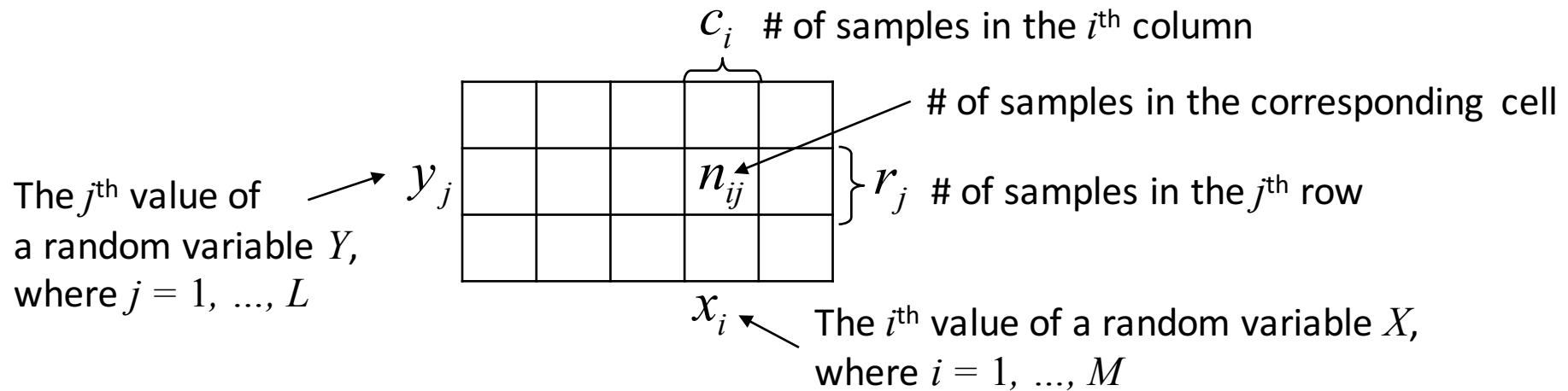
Rules of Probability (1)

- Assume two random variables, X and Y
 - $X : x_1 = \text{“Bread”}, x_2 = \text{“Rice”}, \text{ or } x_3 = \text{“Noodle”}$
 - $Y : y_1 = \text{“Home” or } y_2 = \text{“Restaurant”}$
- Assume the following data samples $\{X, Y\}$:
 $\{\text{Bread, Home}\}, \{\text{Rice, Restaurant}\}, \{\text{Noodle, Home}\}, \{\text{Bread, Restaurant}\},$
 $\{\text{Rice, Restaurant}\}, \{\text{Noodle, Home}\}, \{\text{Bread, Home}\}, \{\text{Rice, Home}\},$
 and $\{\text{Bread, Home}\}$
- Make the following table showing the number of samples

		Number of samples		
	Home	2	1	2
	Restaurant	1	2	0
		Bread	Rice	Noodle

of samples of {Noodle, Home} →

Rules of Probability (2)



Joint probability : $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$

Marginal probability : $p(X = x_i) = \frac{c_i}{N} = \frac{\sum_{j=1}^L n_{ij}}{N}$

Sum rule of probability : $p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j)$

Conditional probability : $p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$

Product rule of probability: $p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N}$
 $= p(Y = y_j | X = x_i)p(X = x_i)$

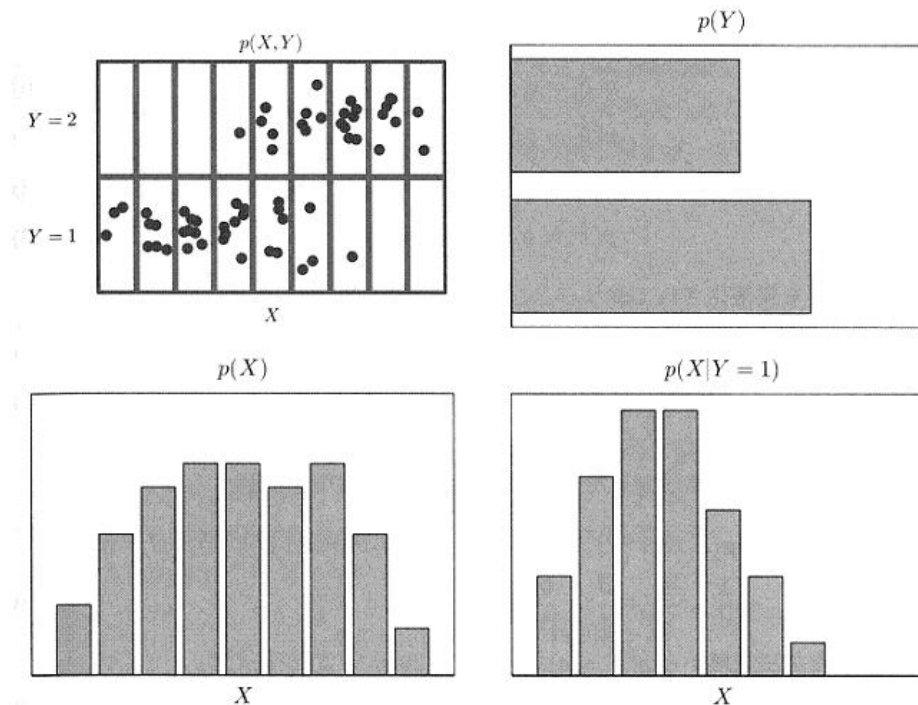
Rules of Probability (3)

- The rules of probability

- **Sum rule** : $p(X) = \sum_Y p(X, Y)$

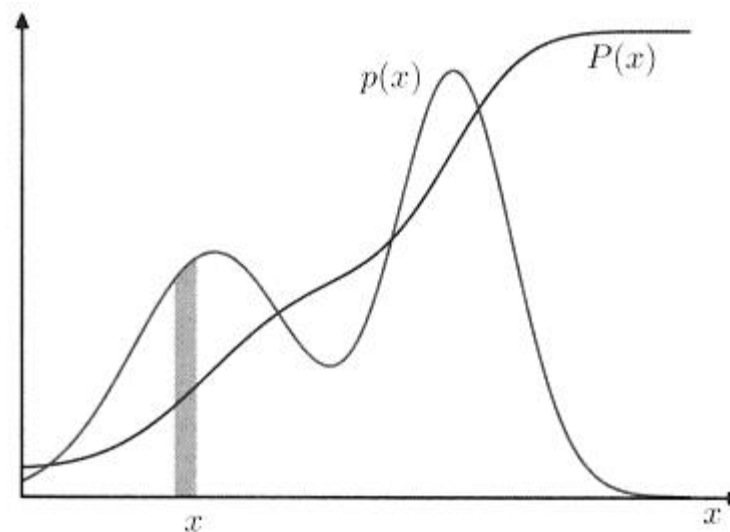
- **Product rule** : $p(X, Y) = p(Y|X)p(X)$

- **Bayes' theorem** : $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$ $\longleftarrow p(X) = \sum_Y p(X|Y)p(Y)$



Probability Densities

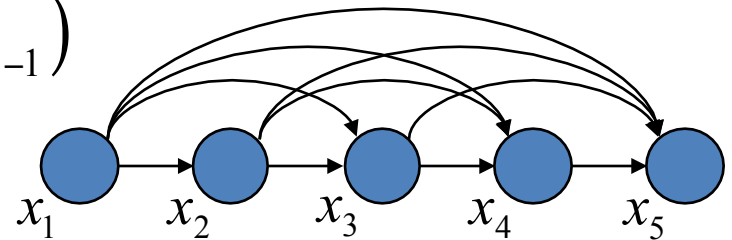
- Probabilities with respect to **continuous variables**
- Probability density over a real-valued variable x : $p(x)$
 - Probability that x will lie in an interval (a, b) : $p(x \in (a, b)) = \int_a^b p(x) dx$
 - Conditions to be satisfied : $\begin{cases} p(x) \geq 0 \\ \int_{-\infty}^{\infty} p(x) dx = 1 \end{cases}$
- Cumulative distribution function: $P(x)$
 - Probability that x lies in the interval $(-\infty, z)$: $P(z) = \int_{-\infty}^z p(x) dx$



How to Model Joint Probability?

- Length of sequential data (# of data points over a sequence) varies...
i.e., # of dimensions of joint probability distribution also varies...
- Joint probability distribution can be represented with conditional probability distributions of individual data points!
i.e., # of distributions varies but # of dimensions of each distribution is fixed.
- However, conditional probability distribution of a present data point given all past data points needs to be modeled...

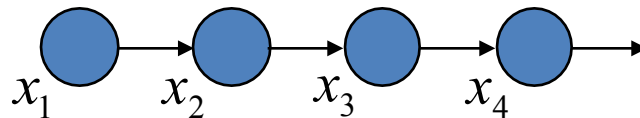
$$\begin{aligned} p(x_1, \dots, x_N) &= p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2) \cdots p(x_N | x_1, \dots, x_{N-1}) \\ &= p(x_1) \prod_{n=2}^N p(x_n | x_1, \dots, x_{n-1}) \end{aligned}$$



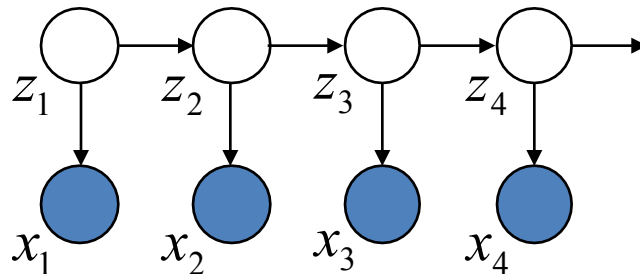
How can we effectively model joint probability distribution of sequential data?

Two Basic Approaches

- **Markov process**



- **Latent variables**



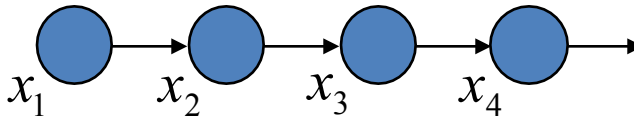
Markov Process

- Assume that the conditional probability distribution of the present states depends only on a few past states

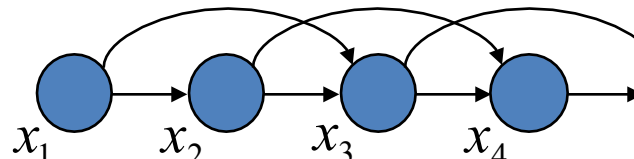
$$p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n | x_1, \dots, x_{n-1})$$

e.g., it depends on
only one past state... $p(x_n | x_1, \dots, x_{n-1}) = p(x_n | x_{n-1})$

1st order Markov chain

$$p(x_1, \dots, x_N) = p(x_1) \prod_{n=2}^N p(x_n | x_{n-1})$$


2nd order Markov chain

$$p(x_1, \dots, x_N) = p(x_1)p(x_2|x_1) \prod_{n=3}^N p(x_n | x_{n-1}, x_{n-2})$$


Example of 1st Order Markov Process

- How many probability distributions are needed if we model English text using the 1st order Markov process?

If only using 27 characters including “space”,

$P(\text{“This sentence is represented by this ...”})$

$= P(T) P(h|T) P(i|h) P(s|i) P(-|s) P(s|-) P(e|s) P(n|e) P(t|n) P(e|t) \dots$

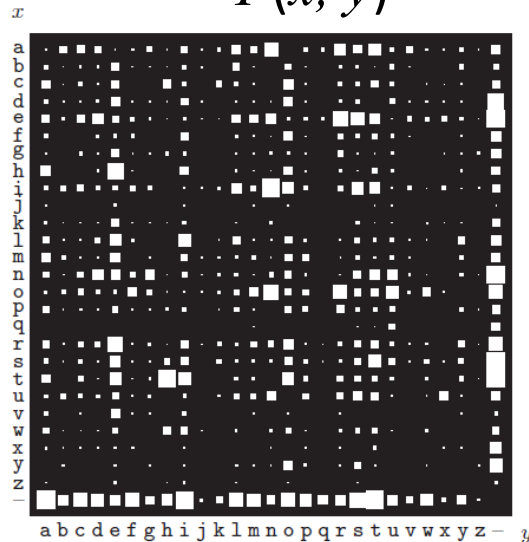
x : 1st letter
 y : 2nd letter

Probability is
shown by the areas
of white squares.

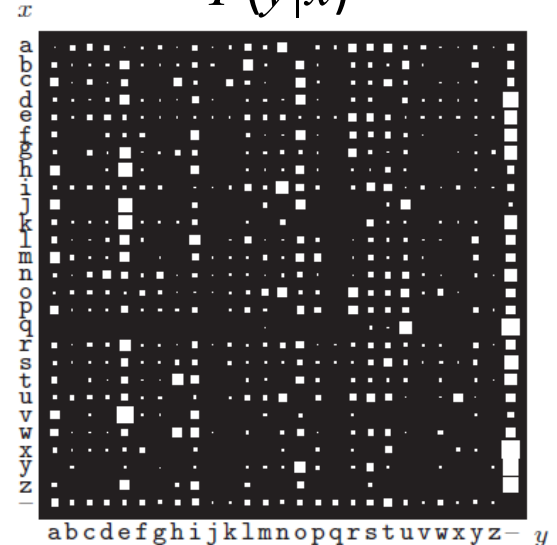
$P(x)$



$P(x, y)$



$P(y|x)$

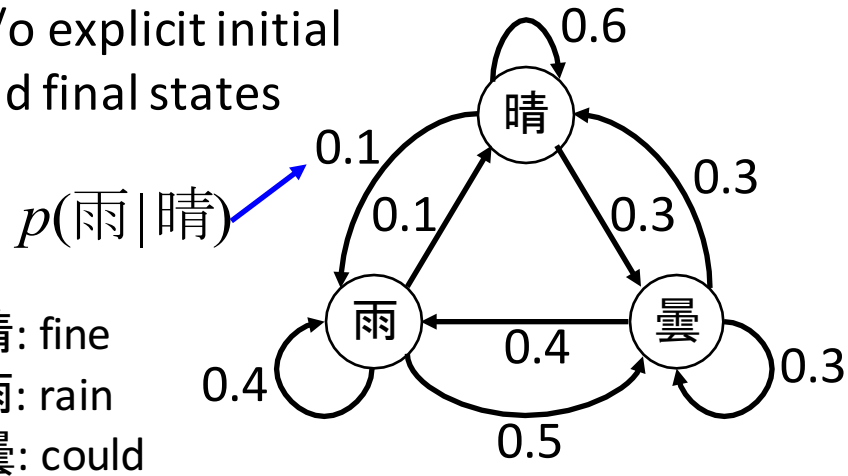


David J.C. MacKay, “Information Theory, Inference, and Learning Algorithms,”
Cambridge University Press, pp. 22 — 24

State Transition Diagram/Matrix

State transition diagram

w/o explicit initial
and final states



$p(\text{雨}|\text{晴})$

State transition matrix

Present state

Past state

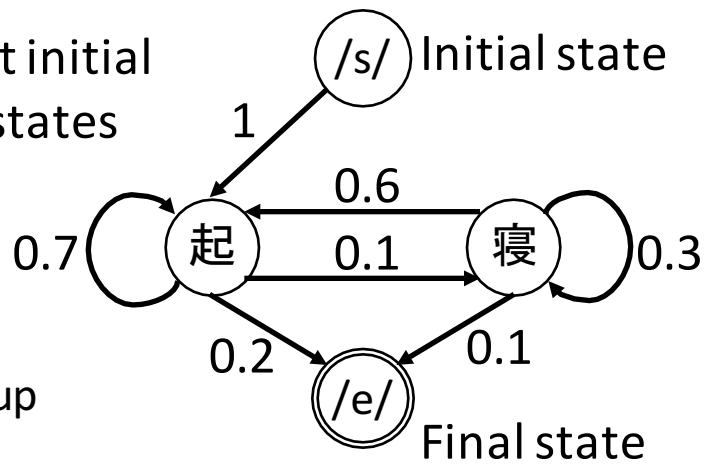
$$A = \begin{matrix} & \begin{matrix} \text{晴} & \text{曇} & \text{雨} \end{matrix} \\ \begin{matrix} \text{晴} \\ \text{曇} \\ \text{雨} \end{matrix} & \begin{bmatrix} 0.6 & 0.3 & \underline{0.1} \\ 0.3 & 0.3 & 0.4 \\ \underline{0.1} & \underline{0.5} & \underline{0.4} \end{bmatrix} \end{matrix}$$

$p(\text{雨}|\text{晴})$

Sum becomes one.

State transition diagram

w/ explicit initial
and final states



State transition matrix

$$A = \begin{matrix} & \begin{matrix} /s/ & 起 & 寢 & /e/ \end{matrix} \\ \begin{matrix} /s/ \\ 起 \\ 寢 \\ /e/ \end{matrix} & \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & & ? & \\ 0 & & & \\ 0 & 0 & 0 & \underline{1} \end{bmatrix} \end{matrix}$$

Final state transition

Example of Applications

- Language model (modeling discrete variables)

- Modeling a word (morpheme) sequence with Markov model (n-gram)

e.g., 「学校に行く」  /s/, 学校, に, 行, く, /e/

Decompose into
morphemes



Modeling with bi-gram (2-gram)

$$p(s)p(\text{学校}|s)p(\text{に}|\text{学校})p(\text{行}|\text{に})p(\text{く}|\text{行})p(e|\text{く})$$

- Autoregressive model (modeling continuous variables)

- Predicting the present data point from past M data points

$$p(x_n | x_{n-M}, \dots, x_{n-1}) = \mathcal{N}(x_n | \underbrace{a_M x_{n-M} + \dots + a_1 x_{n-1}}_{\text{Linear combination of past } M \text{ data points}}, \sigma^2)$$

Linear combination of
past M data points

Model Training (Maximum Likelihood Estimation)

- Training of conditional probability distribution from sequential data samples given as training data $\{x_1^{(1)}, \dots, x_{N_1}^{(1)}\} \dots \{x_1^{(S)}, \dots, x_{N_S}^{(S)}\}$

Likelihood function: $\prod_{s=1}^S p(x_1^{(s)}, \dots, x_{N_s}^{(s)} | \lambda)$

Function of model parameters λ Model parameter set

Determine the conditional probability distributions $p(w_c | w_p, \lambda)$ that maximizes the (log-scaled) likelihood function

$$\arg \max \sum_{s=1}^S \log p(x_1^{(s)} | \lambda) + \sum_{s=1}^S \sum_{n=2}^{N_s} \log p(x_n^{(s)} | x_{n-1}^{(s)}, \lambda)$$

subject to $\sum_{x \in \text{all } w} p(x | w_p, \lambda) = 1$ Constraint to normalize the estimates as probability

Maximum likelihood estimate:

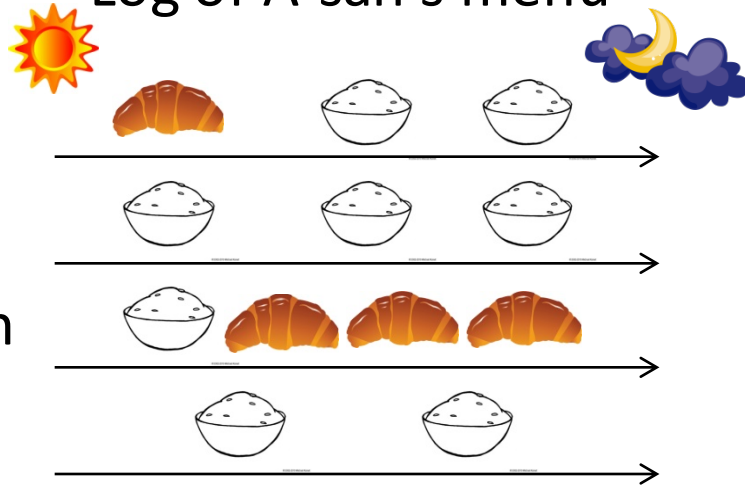
$$p(w_c | w_p, \lambda) = \frac{\text{Num}(w_p, w_c)}{\text{Num}(w_p)}$$

of samples $\{w_p, w_c\}$ # of samples w_p

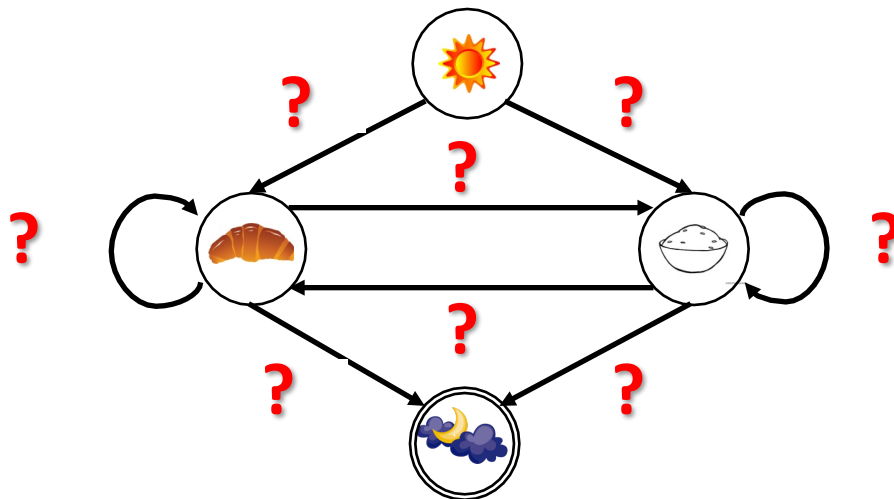
Example of MLE

morning : initial stater

Log of A-san's menu



A-san



MLE of conditional probabilities

$$P(\text{Bread} | \text{Sun}) = ? \quad \text{Count}(\text{Morning, Bread}) / \text{Count}(\text{Morning})$$

$$P(\text{Rice} | \text{Sun}) = ? \quad \text{bCount}(\text{Morning, Rice}) / \text{Count}(\text{Morning})$$

$$P(\text{Bread} | \text{Bread}) = ?$$

$$P(\text{Rice} | \text{Bread}) = ?$$

$$P(\text{Moon/Beans} | \text{Bread}) = ?$$

$$P(\text{Bread} | \text{Rice}) = ?$$

$$P(\text{Rice} | \text{Rice}) = ?$$

$$P(\text{Moon/Beans} | \text{Rice}) = ?$$

State transition diagram

Methods for Evaluating Models

- Use of a test data set $\{w_1, \dots, w_N\}$ not included in training data
- Evaluation metrics

- Likelihood

$$p(w_1, \dots, w_N | \lambda) = \prod_{n=1}^N p(w_n | w_{n-1})$$

- Log-scaled likelihood

$$\log_2 p(w_1, \dots, w_N | \lambda) = \sum_{n=1}^N \log_2 p(w_n | w_{n-1})$$

- Entropy

$$H = -\frac{1}{N} \log_2 p(w_1, \dots, w_N | \lambda) = -\frac{1}{N} \sum_{n=1}^N \log_2 p(w_n | w_{n-1})$$

- Perplexity

$$PP = 2^H$$

A measure of effective “branching factor”

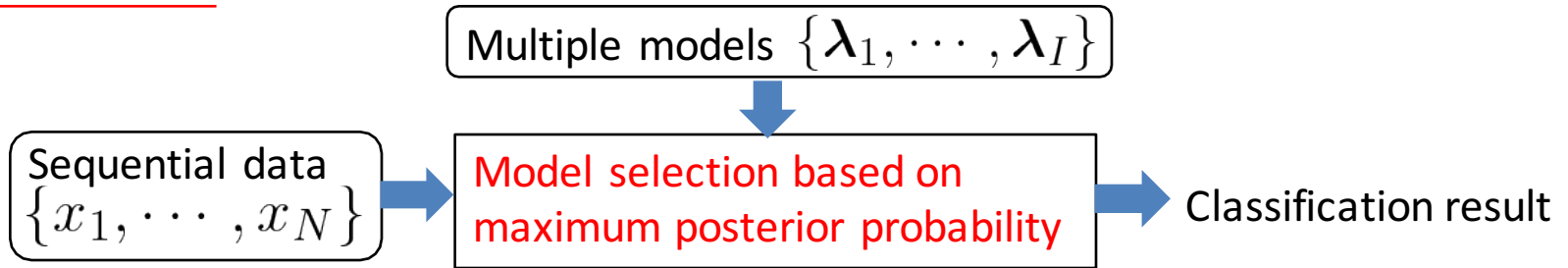
e.g., set a uniform distribution to all n -gram probabilities for M words

$$H = -\frac{1}{N} \sum_{n=1}^N \log_2 \frac{1}{M} = \frac{1}{N} \log_2 M^N = \log_2 M$$

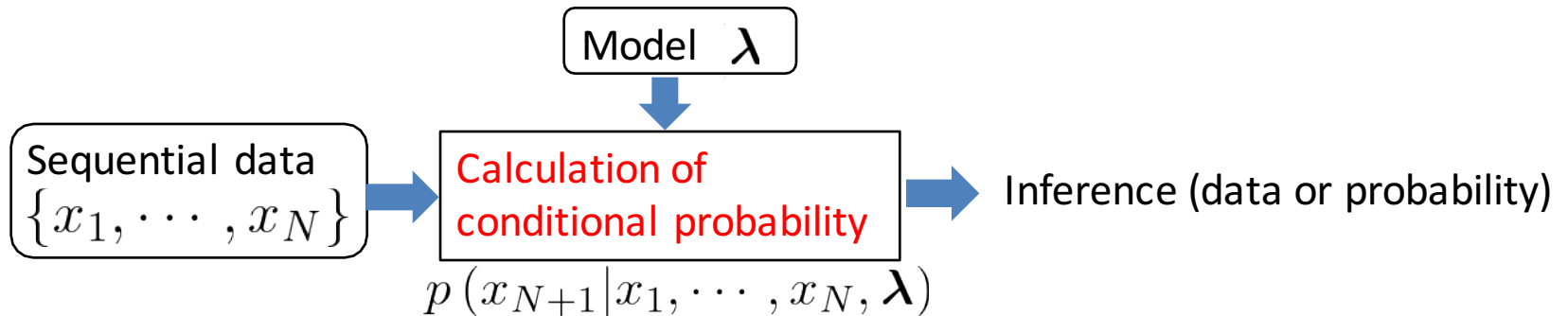
$$PP = M$$

Classification/Inference/Generation

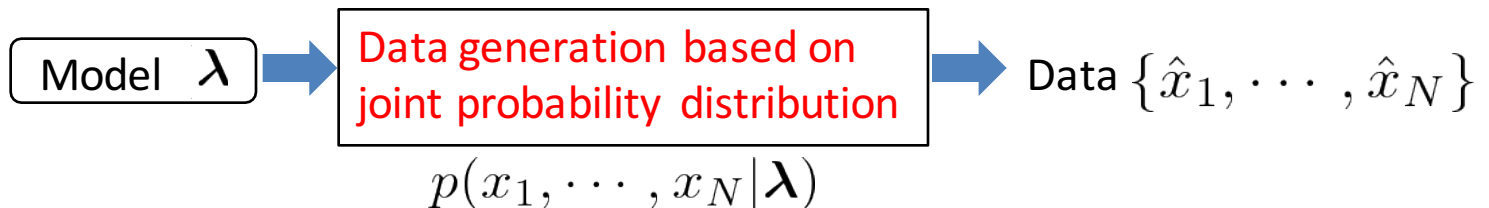
Classification



Inference (prediction)



Generation



Classification w/ Maximum *A Posteriori*

- Select a model that maximizes the posterior probability

Posterior probability:

$$p(\boldsymbol{\lambda}_i | x_1, \dots, x_N) = \frac{\overset{\text{Likelihood function}}{p(x_1, \dots, x_N | \boldsymbol{\lambda}_i)} \overset{\text{Prior probability}}{p(\boldsymbol{\lambda}_i)}}{\underset{\text{Constant}}{p(x_1, \dots, x_N)}}$$

The model is selected by

$$\begin{aligned}\hat{i} &= \arg \max_i p(\lambda_i | x_1, \dots, x_N) \\ &= \arg \max_i p(x_1, \dots, x_N | \lambda_i) p(\lambda_i)\end{aligned}$$

If prior probability is given by a uniform distribution,

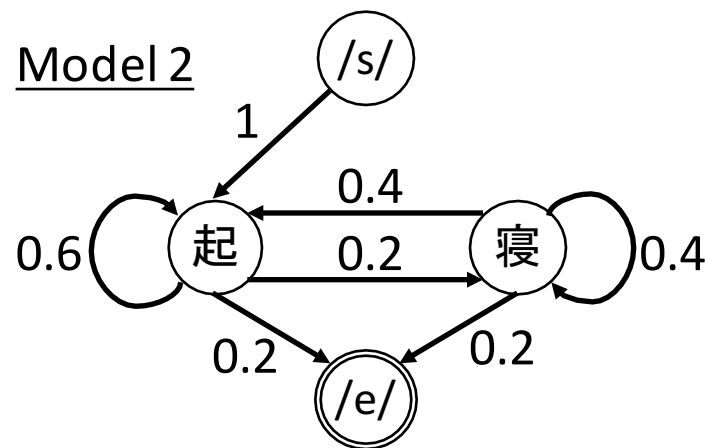
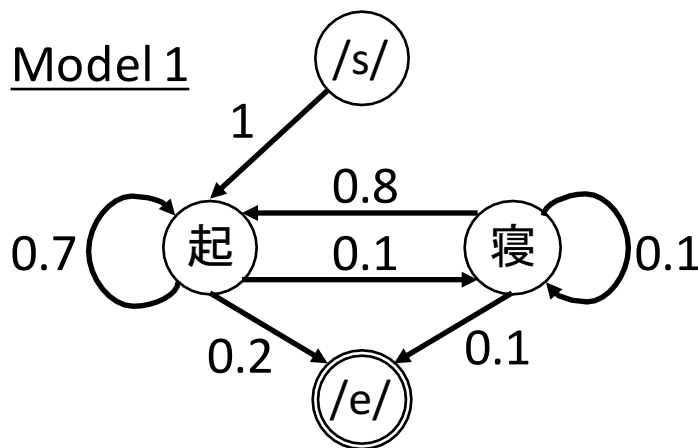
$$\hat{i} = \arg \max_i p(x_1, \dots, x_N | \lambda_i)$$

Classification

- Comparison of model likelihoods

Models:

起: wake up
寝: sleep

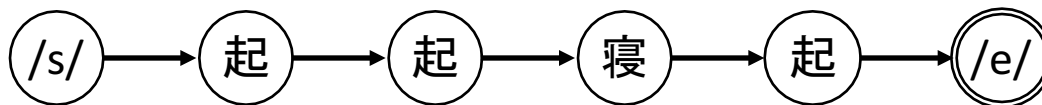


Observed data: /s/ 起 起 寝 起 /e/ Q. Which model is this data sample classified to?

Likelihood: $p(s)p(\text{起}|s)p(\text{起}|\text{起})p(\text{寝}|\text{起})p(\text{起}|\text{寝})p(e|\text{起})$

Trellis:

Model 1: $P(s)=1P(s)=1P(l)=0.7$



Expansion of the state transition graph over time axis

Model 1: $1 \times 1 \times 0.7 \times 0.1 \times 0.8 \times 0.2 = 0.0112$

Model 2: ?

A. Classified to the model 1.

Marginalization for Unobserved Data

- Likelihood calculation with marginalization even if a part of data is not observed.

Observed data: /s/ ? 寝 ? /e/

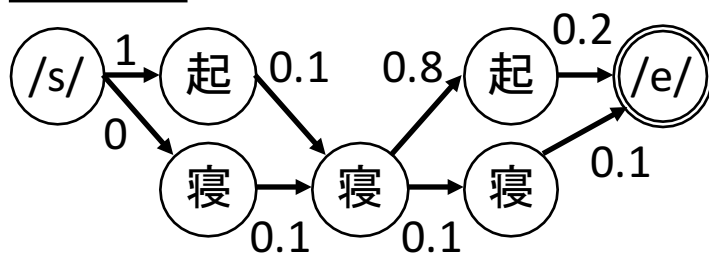
Q. Which model is this data sample classified to?

Likelihood: $\sum_{x_1, x_3 \in \{\text{起}, \text{寝}\}} p(s)p(x_1 | s)p(\text{寝} | x_1)p(x_3 | \text{寝})p(e | x_3)$

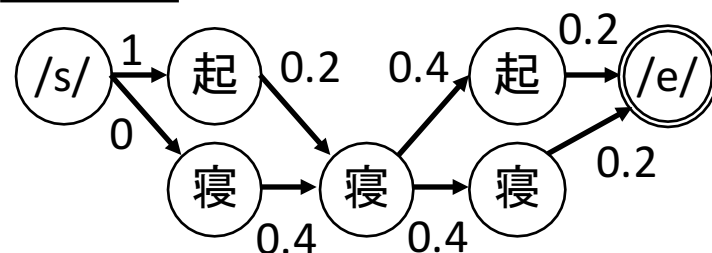
Consider all possible data samples

$$= p(s) \left\{ \sum_{x_1 \in \{\text{起}, \text{寝}\}} p(x_1 | s)p(\text{寝} | x_1) \right\} \left\{ \sum_{x_3 \in \{\text{起}, \text{寝}\}} p(x_3 | \text{寝})p(e | x_3) \right\}$$

Trellis: Model 1



Model 2



$$1 \times (1 \times 0.1 + 0 \times 0.1) \\ \times (0.8 \times 0.2 + 0.1 \times 0.1) = 0.017$$

?

A. Classified to the model 2.

Inference

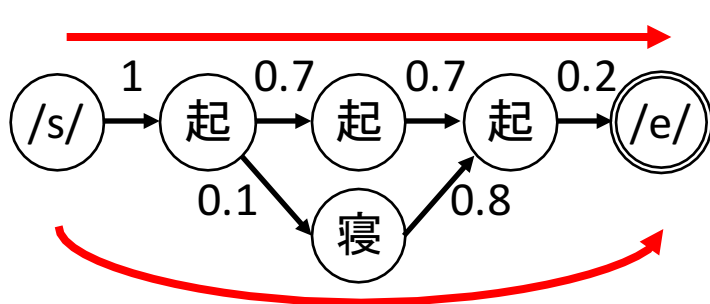
- Calculation of posterior probability

Observed data: /s/ 起 ? 起 /e/

Q. Which “起” or “寢” is likely observed at “?” point?

Posterior probability:

$$p(x_2 = \text{起} | x_1 = \text{起}, x_3 = \text{起}) = \frac{p(s, \text{起}, \text{起}, \text{起}, e)}{\sum_{x_2 \in \{\text{起}, \text{寢}\}} p(s, \text{起}, x_2, \text{起}, e)}$$

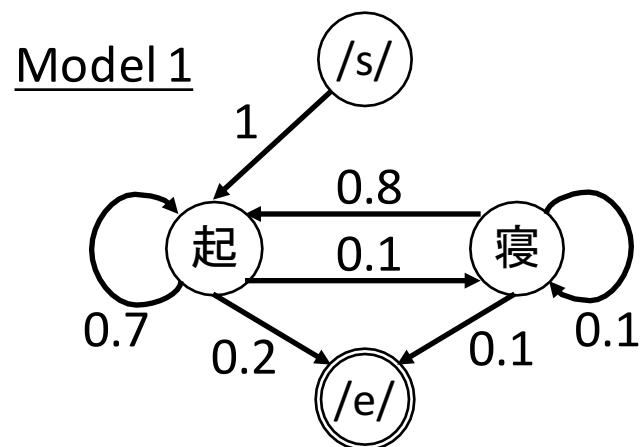


$$p(s, \text{起}, \text{起}, \text{起}, e) = 1 \times 1 \times 0.7 \times 0.7 \times 0.2 = 0.098$$

$$p(s, \text{起}, \text{寢}, \text{起}, e) = ?$$

$$p(x_2 = \text{起} | x_1 = \text{起}, x_3 = \text{起}) = \frac{0.098}{0.098 + 0.016} \approx 0.860$$

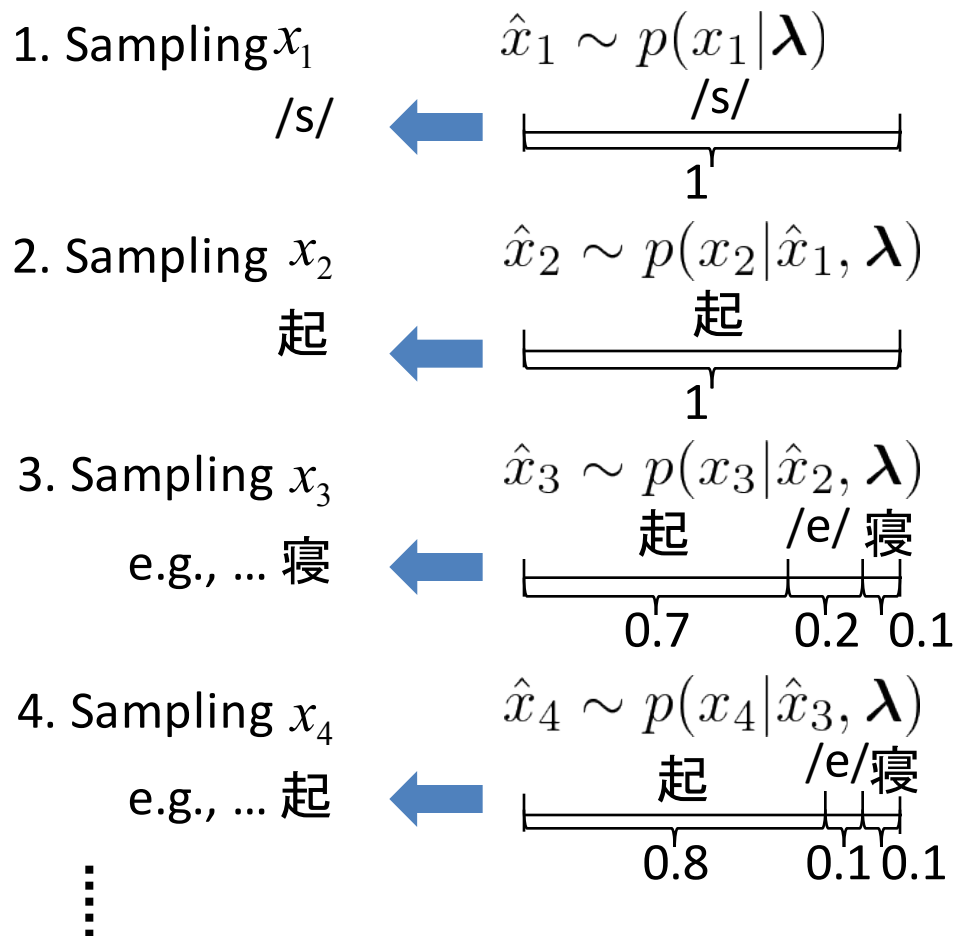
$$p(x_2 = \text{寢} | x_1 = \text{起}, x_3 = \text{起}) = ?$$



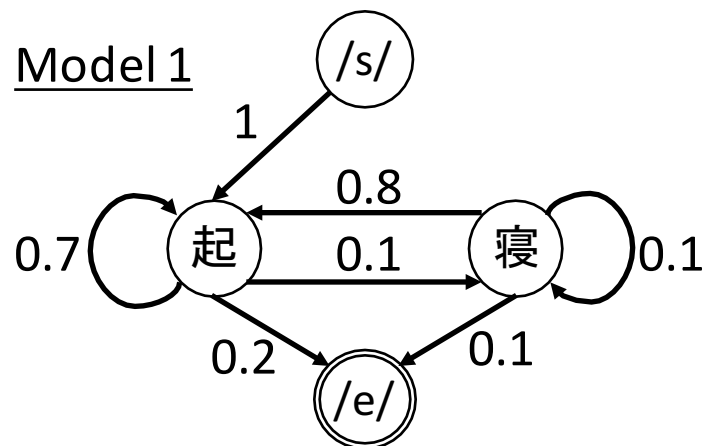
A. “?” is “起” with 86% of probability.

Generation

- Random generation of data samples from the model



End if the final state /e/ is sampled



/s/ 起 寝 起 ... /e/

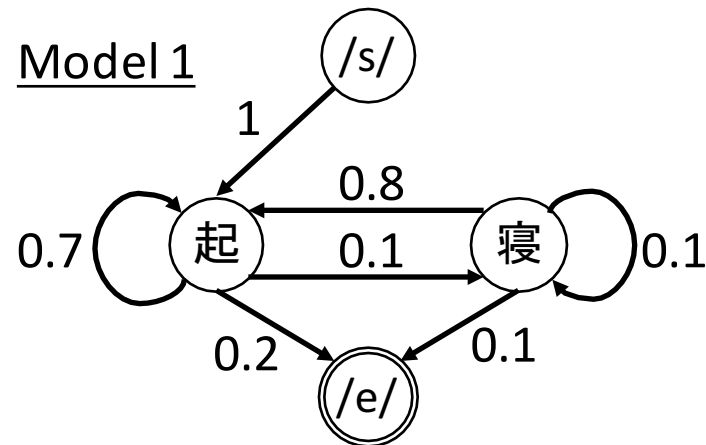
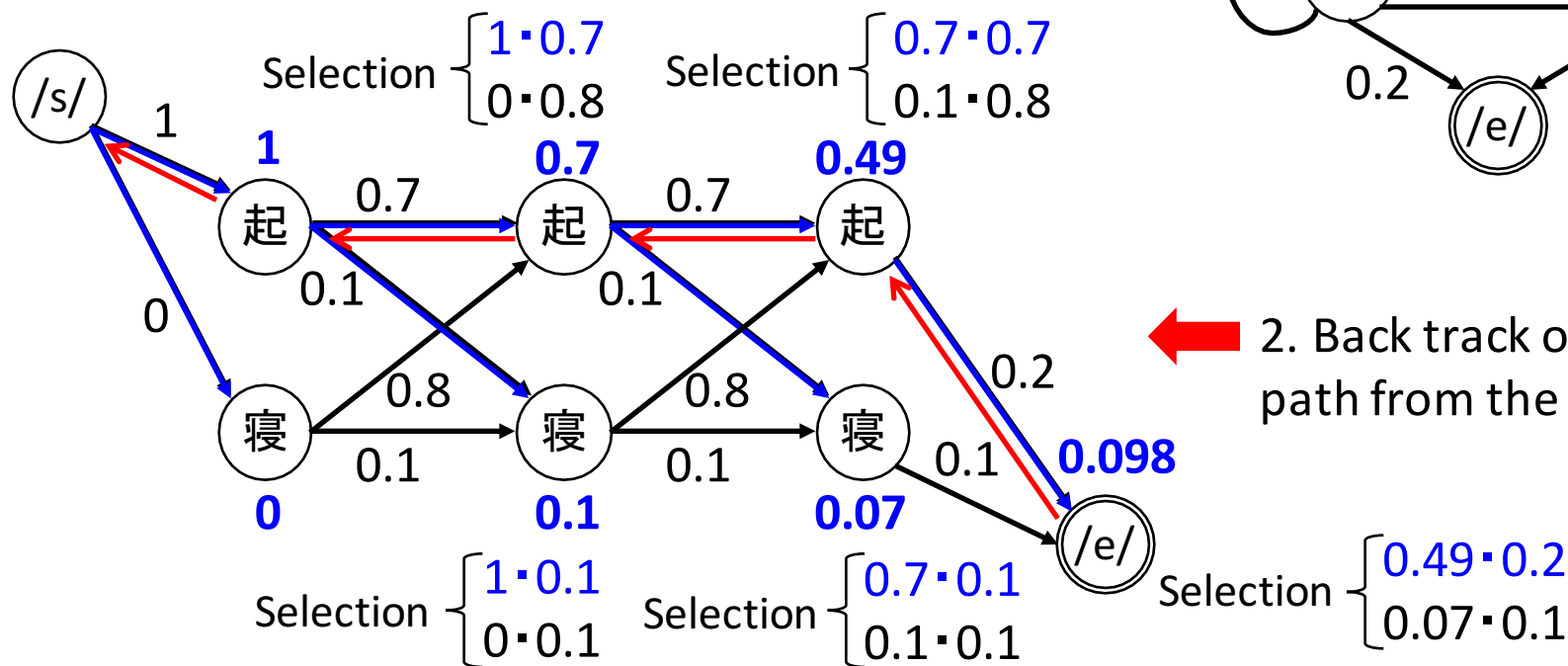
Various lengths of data are generated.

Maximum Likelihood Data Generation

- Data generation by maximizing likelihood under the condition that the length of data is given

Dynamic programming

➡ 1. Store the best path at each state

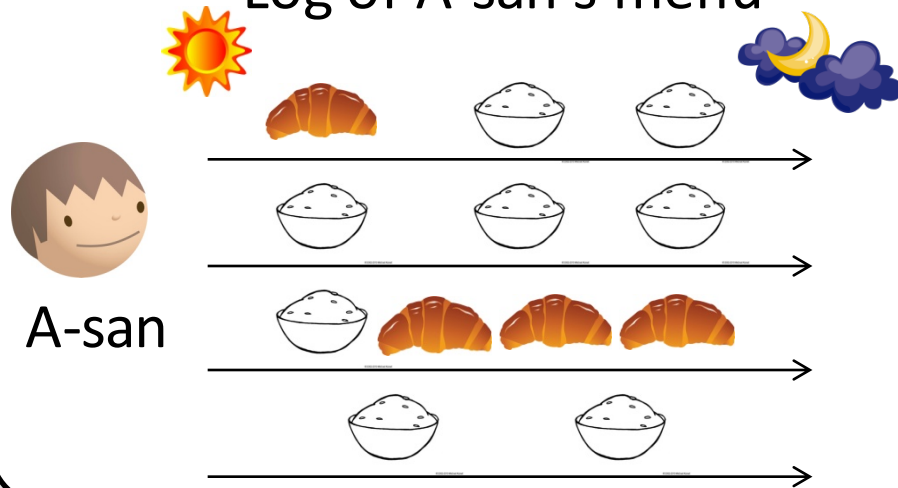


➡ 2. Back track of the best path from the final state

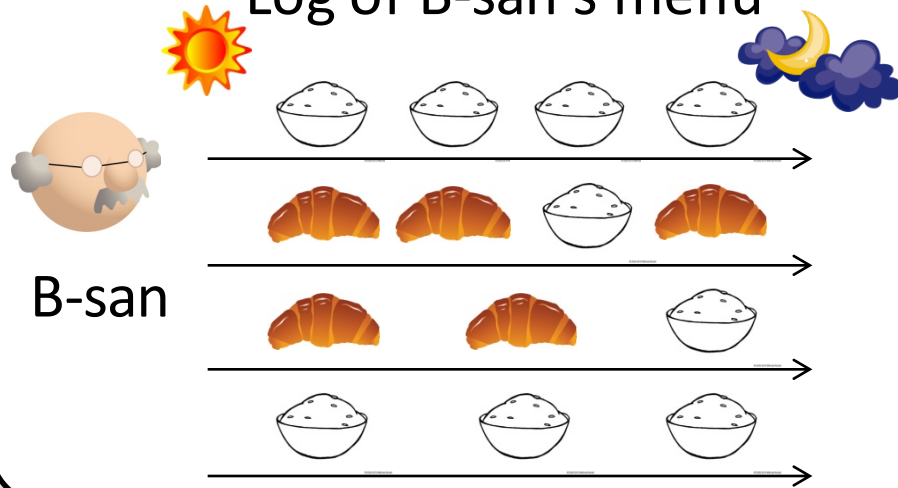
/s/ 起 起 起 /e/ is generated if setting the data length to 3.

You Can Answer the Question!

Log of A-san's menu





Log of B-san's menu



One day, someone ate the following menu.



Q1. A-san or B-san?

Q2. If this is A-san's menu, which is "?",  

Q3. ...
⋮

Let's use Markov model to answer them!