

# Probability and Time: Markov Models

Computer Science cpsc322, Lecture 31  
*(Textbook Chpt 6.5.1)*

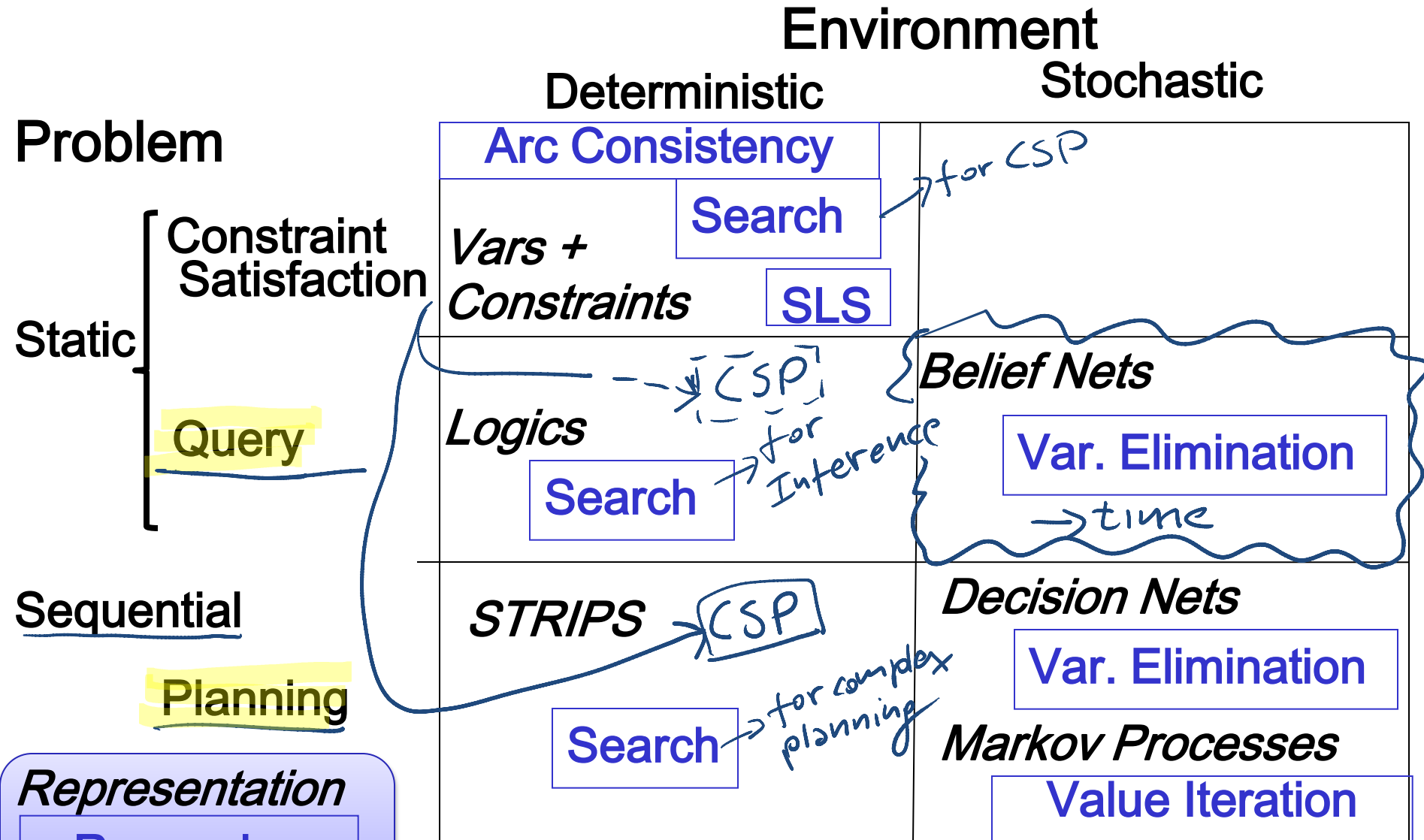
Nov, 22, 2013



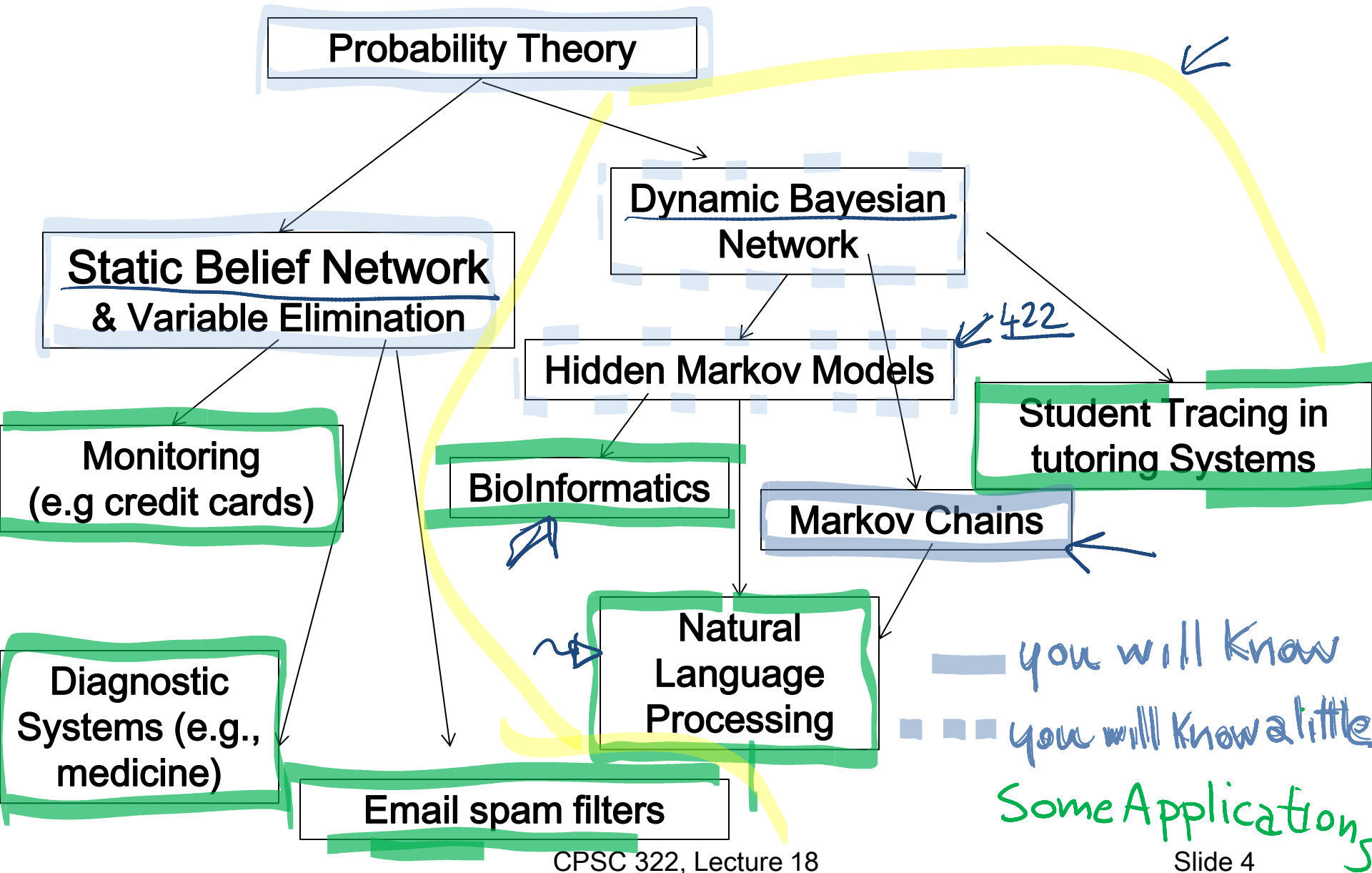
# Lecture Overview

- **Recap**
- Temporal Probabilistic Models
- Start Markov Models
  - Markov Chain
  - Markov Chains in Natural Language Processing

# Big Picture: R&R systems



# Answering Query under Uncertainty



# Lecture Overview

- Recap
- **Temporal Probabilistic Models**
- Start Markov Models
  - Markov Chain
  - Markov Chains in Natural Language Processing

# Modelling static Environments

So far we have used Bnets to perform inference in static environments

- For instance, the system keeps collecting evidence to diagnose the cause of a fault in a system (e.g., a car).



- The environment (values of the evidence, the true cause) does not change as I gather new evidence

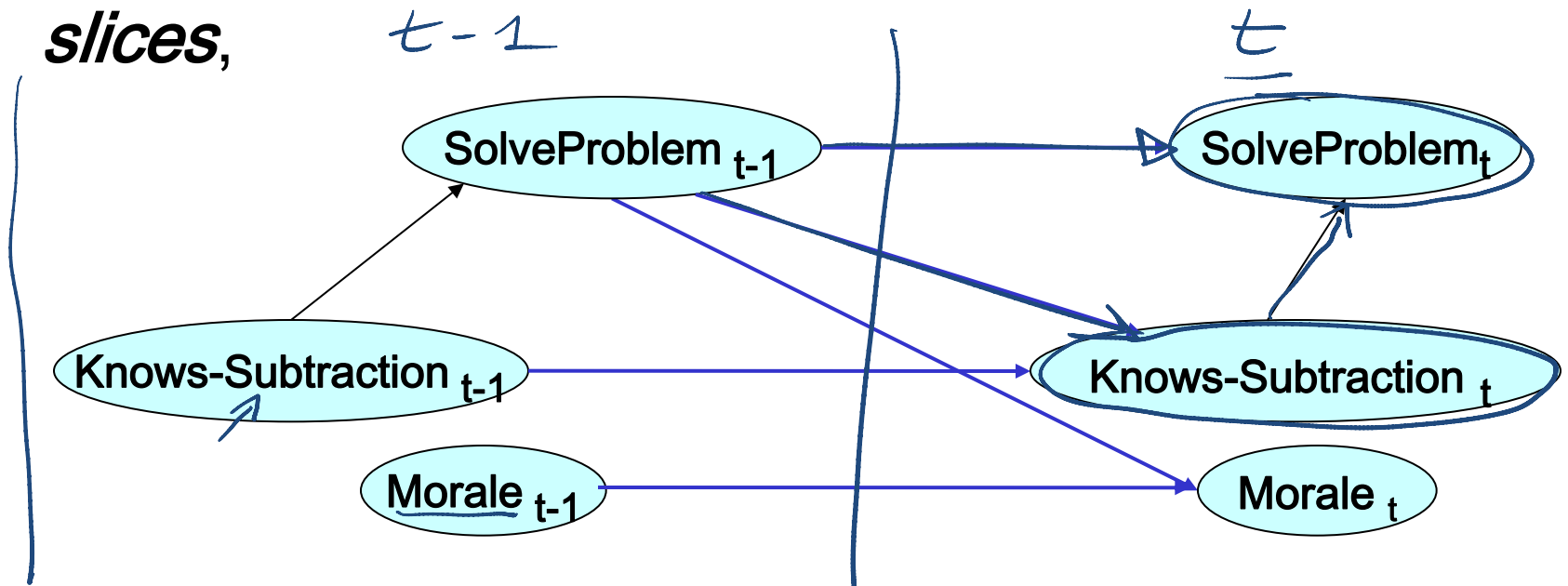
- What does change?

*The system's beliefs over possible causes*



# Modeling Evolving Environments

- Often we need to make inferences about evolving environments.
- Represent the state of the world at each specific point in time via a series of snapshots, or *time slices*,



Tutoring system tracing student *knowledge* and *morale*

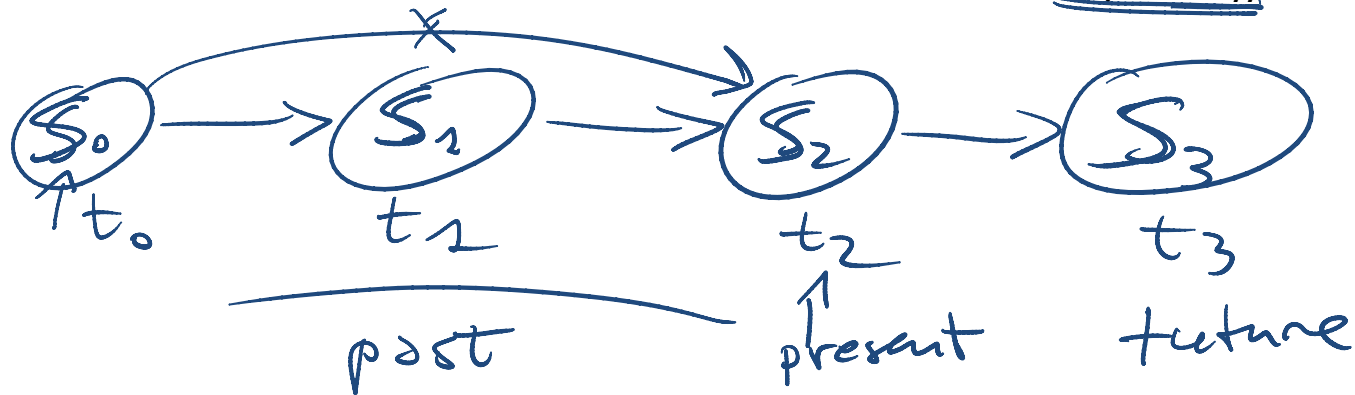
# Lecture Overview

- Recap
- Temporal Probabilistic Models
- Start Markov Models
  - **Markov Chain**
  - Markov Chains in Natural Language Processing



# Simplest Possible DBN

- One random variable for each time slice: let's assume  $S_t$  represents the **state** at time  $t$  with domain  $\{V_1 \dots V_n\}$

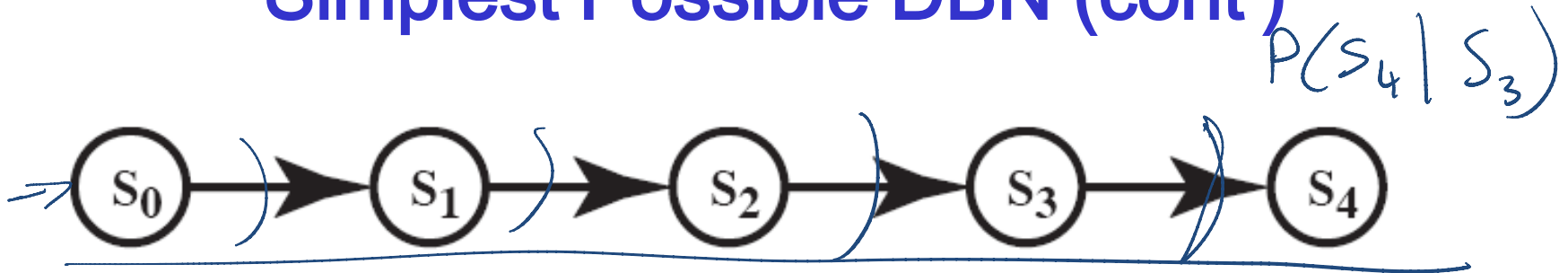


- Each random variable depends only on the previous one

- Thus 
$$P(S_{t+1} | S_0 \dots S_t) = P(S_{t+1} | \underline{S_t})$$

- Intuitively  $S_t$  conveys all of the information about the history that can affect the future states.
- → “The future is independent of the past given the present.”

# Simplest Possible DBN (cont')



- How many CPTs do we need to specify?

iclicker.

4  $P(S_1 | S_0)$   $P(S_2 | S_1)$  etc.

A. 1

C. 2

D. 3

B. 4

- *Stationary process assumption*: the mechanism that regulates how state variables change overtime is **stationary**, that is it can be described by a single transition model
- $P(S_t | S_{t-1})$  is the same for all  $t$

# Stationary Markov Chain (SMC)



A stationary Markov Chain : for all  $t > 0$

- $P(S_{t+1} | S_0, \dots, S_t) = P(S_{t+1} | S_t)$  and *Markov assumption*
- $P(S_{t+1} | S_t)$  is the same *stationary*

We only need to specify  $P(S_0)$  and  $P(S_{t+1} | S_t)$

- Simple Model, easy to specify ←
- Often the natural model ←
- The network can extend indefinitely ←
- **Variations of SMC are at the core of most Natural Language Processing (NLP) applications!** *also used in the PageRank algo (used by Google to rank web pages)*

# Stationary Markov Chain (SMC)



A stationary Markov Chain : for all  $t > 0$

- $P(S_{t+1} | S_0, \dots, S_t) = P(S_{t+1} | S_t)$  and *Markov assumption*
- $P(S_{t+1} | S_t)$  is the same *stationary*

So we only need to specify?

iclicker.

**A.**  $P(S_{t+1} | S_t)$  and  $P(S_0)$

**B.**  $P(S_0)$

**C.**  $P(S_{t+1} | S_t)$

**D.**  $P(S_t | S_{t+1})$

# Stationary Markov-Chain: Example

Domain of variable  $S_i$  is  $\{t, q, p, a, h, e\}$

Probability of initial state  $P(S_0)$

Stochastic Transition Matrix  $P(S_{t+1}|S_t)$

Which of these two is a possible STM?

t	.6
q	.4
p	0
a	0
h	0
e	

	$S_{t+1}$					
	t	q	p	a	h	e
$S_t$ t	0	.3	0	.3	.4	0
q	.4	0	.6	0	0	0
p	0	0	1	0	0	0
a	0	0	.4	.6	0	0
h	0	0	0	0	0	1
e	1	0	0	0	0	0

	$S_{t+1}$					
	t	q	p	a	h	e
$S_t$ t	1	0	0	0	0	0
q	0	1	0	0	0	0
p	.3	0	1	0	0	0
a	0	0	0	1	0	0
h	0	0	0	0	0	1
e	0	0	0	.2	0	1

⊖  
Σ > 1

**A.** Left one only

**C.** Both

**B.** Right one only

**D.** None

# Stationary Markov-Chain: Example

Domain of variable  $S_i$  is  $\{t, q, p, a, h, e\}$

six possible values

We only need to specify...

$$P(S_0)$$

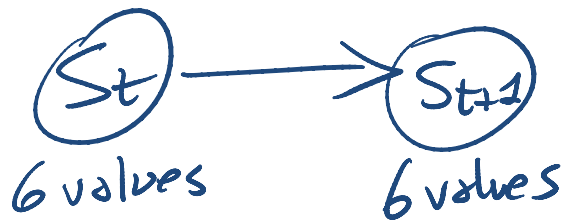
Probability of initial state

t	.6
q	.4
p	0
a	0
h	0
e	0

Stochastic Transition Matrix

$$P(S_{t+1}|S_t)$$

$S_{t+1}$



	t	q	p	a	h	e
t	0	.3	0	.3	.4	0
q	.4	0	.6	0	0	0
p	0	0	1	0	0	0
a	0	0	.4	.6	0	0
h	0	0	0	0	0	1
e	1	0	0	0	0	0

$\leftarrow P(S_{t+1}|S_t=q)$   
 $\leftarrow P(S_{t+1}|S_t=p)$

...

# Markov-Chain: Inference

Probability of a sequence of states  $S_0 \dots S_T$

$$P(S_0, \dots, S_T) = P(S_0) P(S_1 | S_0) P(S_2 | S_1) \dots$$



$P(\text{u, e, e})$

$P(S_0)$

t	.6
q	.4
p	0
a	0
h	0
e	0

$P(S_{t+1} | S_t)$

	t	q	p	a	h	e
t	0	.3	0	.3	.4	0
q	.4	0	.6	0	0	0
p	0	0	1	0	0	0
a	0	0	.4	.6	0	0
h	0	0	0	0	0	1
e	1	0	0	0	0	0

**Example:**

$$P(t, q, p) =$$

$$P(t) * P(q|t) * P(p|q) = .6 * .3 * .6 = .108$$

# Lecture Overview

- Recap
- Temporal Probabilistic Models
- **Markov Models**
  - Markov Chain
  - Markov Chains in Natural Language Processing



# Key problems in NLP

Noun Verb

*"Book me a room near UBC"*

$w_1$   $w_2$   $w_3$   $w_4$   $w_5$   $w_6$

$$P(w_1, \dots, w_n)?$$

Assign a probability to a sentence (a sequence of words)

- • Part-of-speech tagging → **Summarization, Machine**
- • Word-sense disambiguation, → **Translation.....**
- Probabilistic Parsing

Predict the next word

$$P(w_n | w_1 \dots w_{n-1}) = \\ = P(w_1 \dots w_n) / P(w_1 \dots w_{n-1})$$

- • Speech recognition
- • Hand-writing recognition
- • Augmentative communication for the disabled

$$P(w_1, \dots, w_n)?$$

**Impossible to estimate ☹**

$P(w_1, \dots, w_n)$ ?

Impossible to estimate!

Assuming  $10^5$  words and average sentence contains 10 words .....

$(10^5)^{10} = 10^{50}$   
would contain  $\uparrow$  probabilities

**Google language repository** (22 Sept. 2006)  $\rightarrow$  collected from the whole web  
contained "only": 95,119,665,584 sentences  
 $\sim 10^{11}$

Most sentences will not appear or appear only once ☹

# What can we do?

Make a strong simplifying assumption!

Sentences are generated by a Markov Chain

$$\begin{aligned} P(w_1, \dots, w_n) &= \overbrace{P(w_1 | \langle S \rangle)}^{w_1 \text{ at the beginning of a sentence}} \prod_{k=2}^n P(w_k | w_{k-1}) \\ &= P(w_1 | \langle S \rangle) P(w_2 | w_1) P(w_3 | w_2) \dots P(w_n | w_{n-1}) \end{aligned}$$

**P(The big red dog barks)=**

$$\begin{aligned} &P(\underline{\text{The}} | \langle S \rangle) * P(\text{big} | \text{the}) * P(\text{red} | \text{big}) * \dots \\ &* P(\text{dog} | \text{red}) * P(\text{barks} | \text{dog}) \end{aligned}$$

These probs can be assessed in practice!



# Estimates for Bigrams

$$P(w_i | w_{i-1})$$

Silly language repositories with only two sentences:

"<S> The big red dog barks against the big pink dog"

"<S> The big pink dog is much smaller"

Count how many times in your documents you have "big red" and "big"

$$P(\underline{red} | \underline{big}) = \frac{P(\underline{big}, \underline{red})}{P(\underline{big})} = \frac{\frac{\overset{\text{count}}{C(\underline{big}, \underline{red})}}{\cancel{N_{\text{pairs}}}}}{\frac{C(\underline{big})}{\cancel{N_{\text{words}}}}} = \frac{C(\underline{big}, \underline{red})}{C(\underline{big})} = \frac{1}{3}$$

$P(w_i | w_{i-1})$   
 $10^5 * 10^5$  matrix

$$P(w_i | w_{i-2}, w_{i-1})$$

some models use two preceding words

# Bigrams in practice...

If you have  $10^5$  words in your dictionary  $P(w_i | w_{i-1})$

will contain this many numbers.... ??

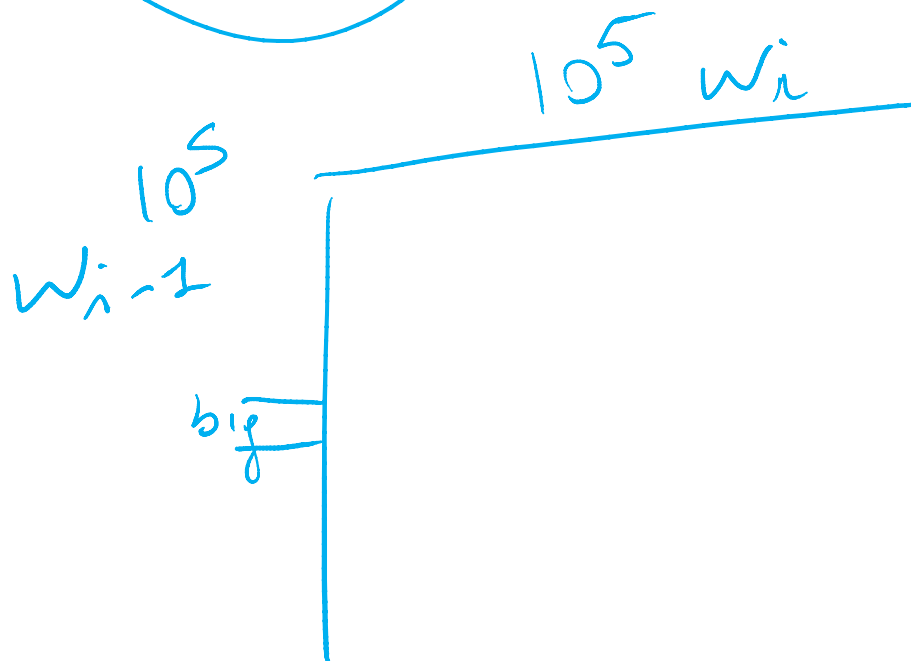
A.  $2 * 10^5$

B.  $10^{10}$

C.  $5 * 10^5$

D.  $2 * 10^{10}$

iclicker.



$$10^5 * 10^5$$

# Learning Goals for today's class

You can:

- Specify a Markov Chain and compute the probability of a sequence of states
- Justify and apply Markov Chains to compute the probability of a Natural Language sentence  
(NOT to compute the conditional probabilities - slide 18)

# Markov Models

Markov Chains

*Simplest Possible  
Dynamic Bnet*

Hidden Markov  
Model

*We cannot observe  
directly what we care  
about*

*Add Actions and  
Values (Rewards)*

422  
Markov Decision  
Processes (MDPs)

## Next Class

- **Finish Probability and Time:** Hidden Markov Models (HMM) (*TextBook 6.5.2*)
- **Start Decision networks** (*TextBook chpt 9*)

## Course Elements

- **Assignment 4 is available on Connect Due on Dec the 2<sup>nd</sup> .**

**Office Hours today**

**2-3 => 3-4**