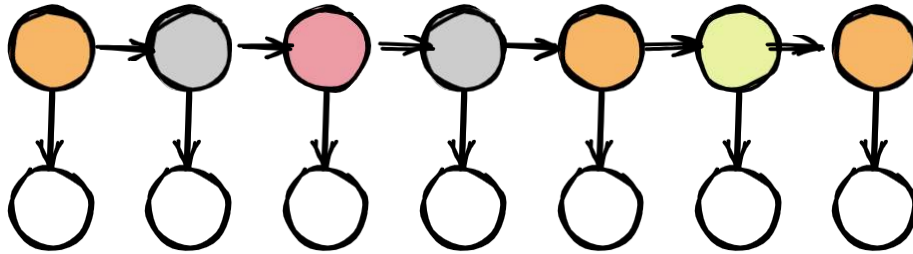# Hidden Markov Models and Sequence Data

Week 18

Middlesex University Dubai; CST4050 Fall21;
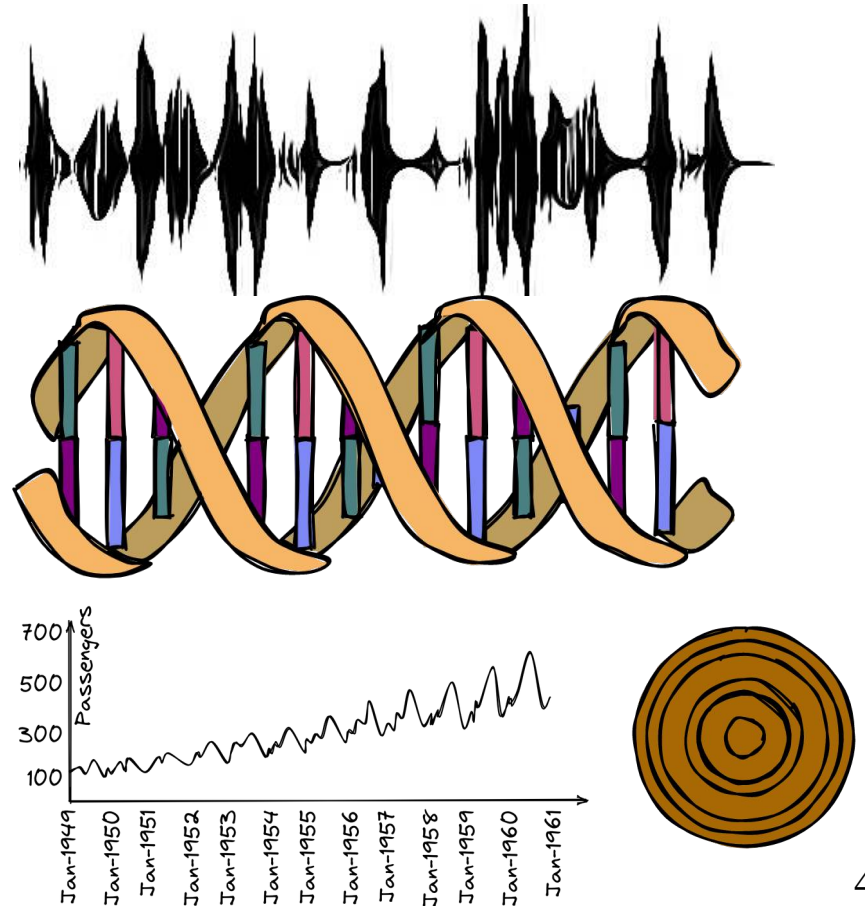Instructor: Dr. Ivan Reznikov

# Plan

- Sequential Data

- Sequential Labeling

- Bayesian Networks

- Mixture Models

- Markov assumption

- Hidden Markov Model

# What is sequence data?

- Ordered set of elements: $x = x_1, x_2, \ldots x_N$

- Order determined by time or position and could be regular or irregular

- Each element $x_i$ could be

  - Numerical (sales, stock price, etc.)
  - Categorical (weather, part-of-speech)
  - Multiple attributes

- The length N of a sequence isn't fixed
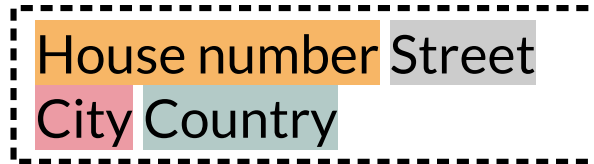
# Examples of sequence data

- Speech (sequence of phonemes)

- Language-related (sequence of words)

- Bioinformatics (genes – sequence of 4 possible nucleotides and proteins – sequence of 20 possible amino-acids)

- Telecommunications (sequence of data packets)

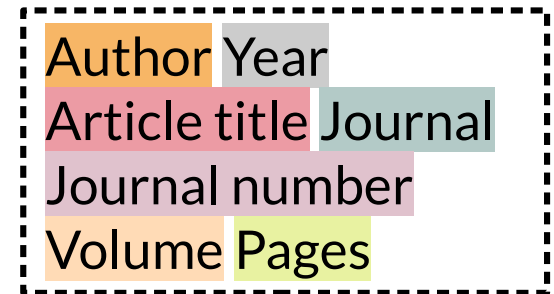- Time series (sequence of events per time)

- …

# Sequence labeling

Address:
221B Baker Street, London, UK

House number Street
City Country

Citation:
Pauling, L. (1931). The nature of the chemical bond. II. The one-electron bond and the three-electron bond. Journal of the American Chemical Society, 53(9), 3225-3237.

Author Year
Article title Journal
Journal number
Volume Pages

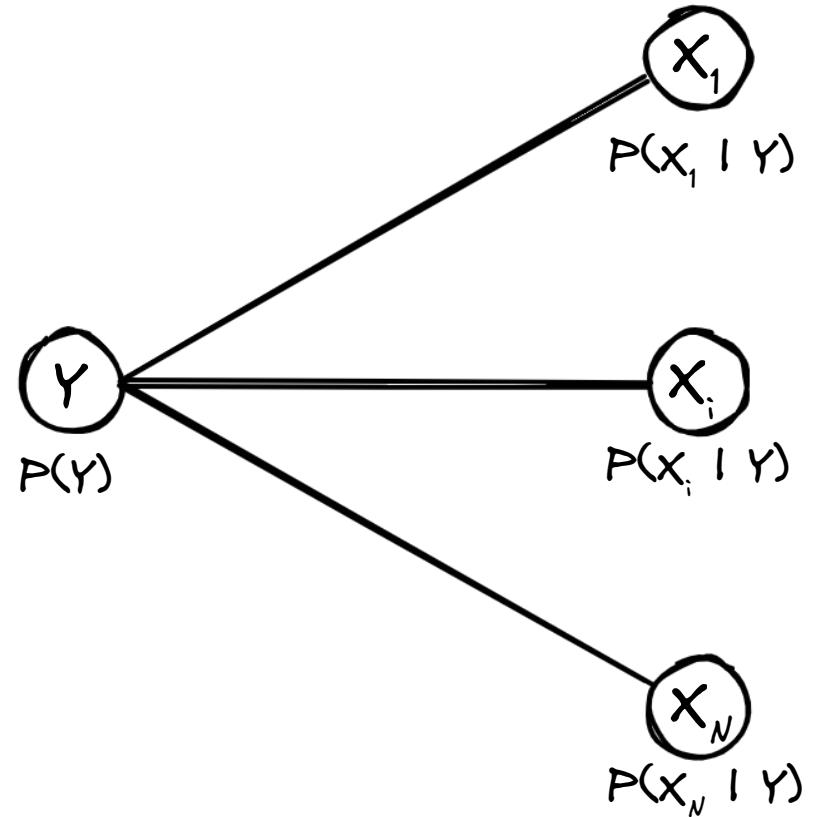**Input**: a sequence $x = (x_1, \ldots x_n)$

**Output**: a sequence $y = (y_1, \ldots y_n)$, where $y_i$ is a label for $x_i$

# Graphical model

Let's assume we have a condition Y. There are several X, that can occur with Y happening. We can draw represent our graph as a probability tree:

- Edges showing dependencies

- Each node has associated conditional

- Probability distribution, conditioned on its parent nodes

- Nodes are independent

$$P(X_1, X_2, ... X_N, Y) = P(X_1 \mid Y) \times P(X_2 \mid Y) ... P(X_N \mid Y) \times P(Y)$$

$X_1$

$P(X_1 \mid Y)$

$Y$

$P(Y)$

$X_i$

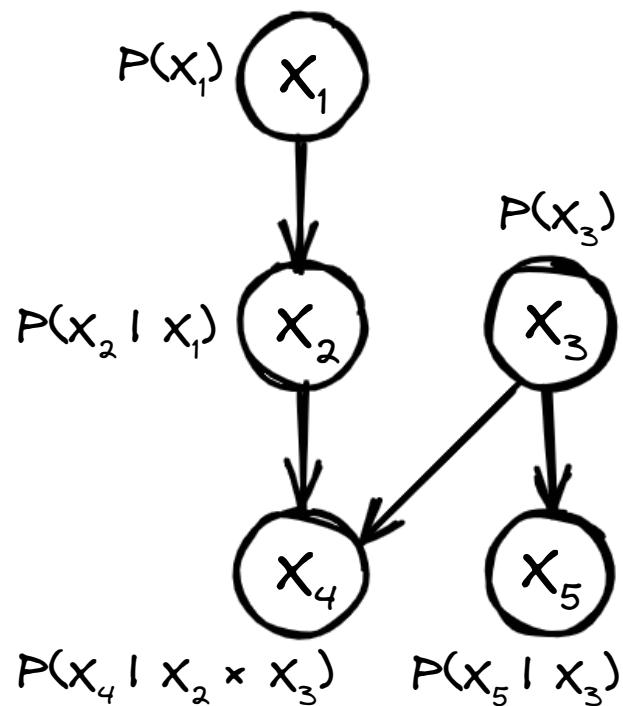$P(X_i \mid Y)$

$X_N$

$P(X_N \mid Y)$

# Graphical model

Let's now draw a directed graph out of 5 nodes.

$$P(X_1, X_2, X_3, X_4, X_5) = P(X_5 \mid X_3) \times$$
$$\times \; P(X_4 \mid X_2 \times X_3) \times P(X_2 \mid X_1) \times$$
$$\times \; P(X_3) \times P(X_1)$$

conditional distributions
marginal distributions



$P(X_1)$   $X_1$

$P(X_3)$

$P(X_2 \mid X_1)$   $X_2$    $X_3$

$X_4$    $X_5$

$P(X_4 \mid X_2 \times X_3)$    $P(X_5 \mid X_3)$

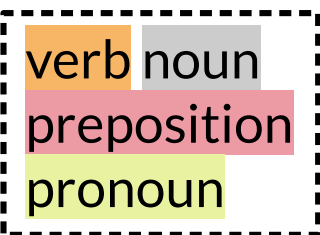# Bayesian Networks

Learning this Bayesian network is equivalent to learning 5 small/simple independent networks from the same data:



$P(x_1)$ $X_1$

$P(x_2 \mid x_1)$ $X_2$

$P(x_3)$ $X_3$

$P(x_4 \mid x_2 \times x_3)$ $X_4$ $X_5$ $P(x_5 \mid x_3)$

8

# Mixture model

whisper words of wisdom let it be

verb noun
preposition
pronoun
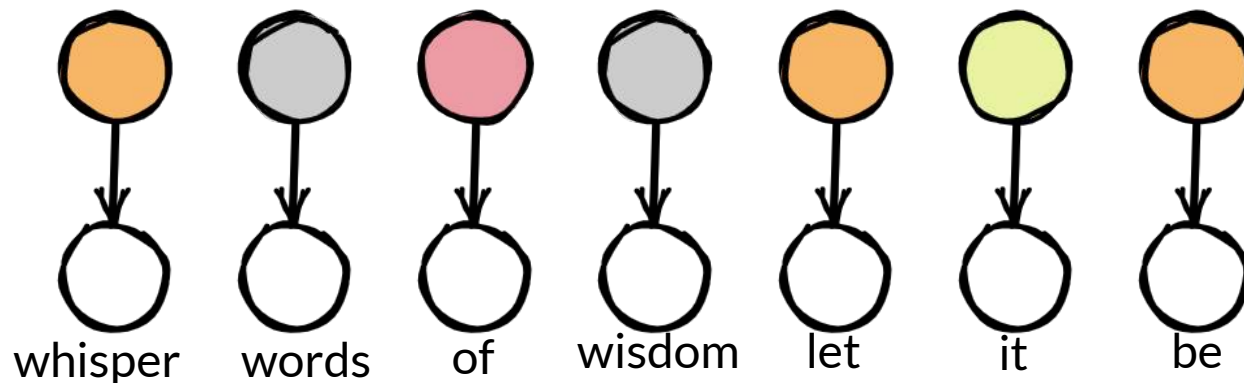
$P(y, x) = P($verb, noun, preposition, noun, verb, pronoun, verb, whisper, words, of, wisdom, let, it, be$) = P($verb, whisper$) \times P($noun, words$) \times$

$\times ...$

$= P($whisper | verb$) \times P($verb$) \times$

$\times P($words | noun$) \times P($noun$) \times$

$\times ...$

whisper    words    of    wisdom    let    it    be

# Mixture model

whisper words of wisdom let it be

|  | whispers | words | of | wisdom | let | it | be |
|---|---|---|---|---|---|---|---|
| verb (0.35) | <u>0.7</u> | 0.2 | 0.1 | 0.05 | <u>0.6</u> | 0.0 | <u>0.9</u> |
| noun (0.4) | 0.3 | <u>0.7</u> | 0.1 | <u>0.85</u> | 0.3 | 0.15 | 0.0 |
| prep (0.15) | 0.0 | 0.0 | <u>0.7</u> | 0.0 | 0.05 | 0.1 | 0.1 |
| pronoun (0.1) | 0.0 | 0.1 | 0.1 | 0.1 | 0.05 | <u>0.65</u> | 0.0 |

whisper   words   of   wisdom   let   it   be

$P(y, x) = P(\text{verb, noun, preposition, noun,}$
$\text{verb, pronoun, verb, whisper, words, of,}$
$\text{wisdom, let, it, be}) = P(\text{verb, whisper}) \times$
$P(\text{noun, words}) \times$
$\times \dots$
$= P(\text{whisper | verb}) \times P(\text{verb}) \times$
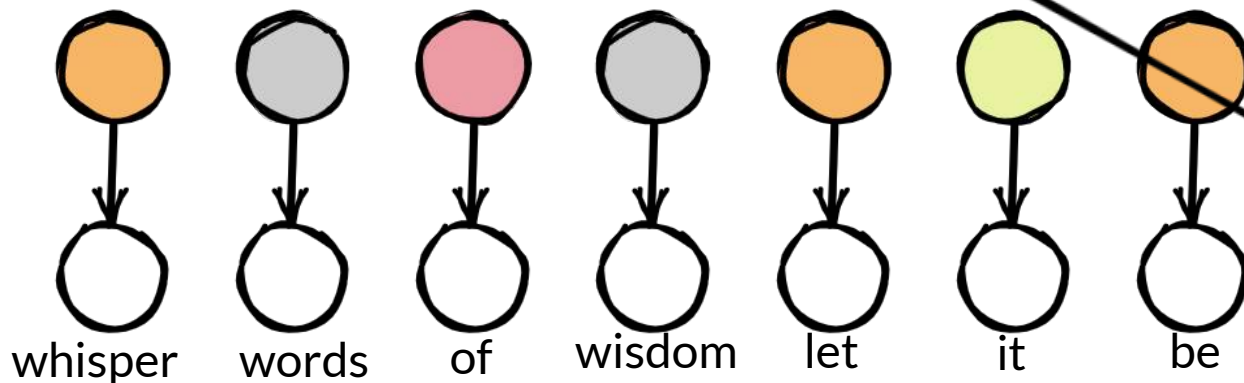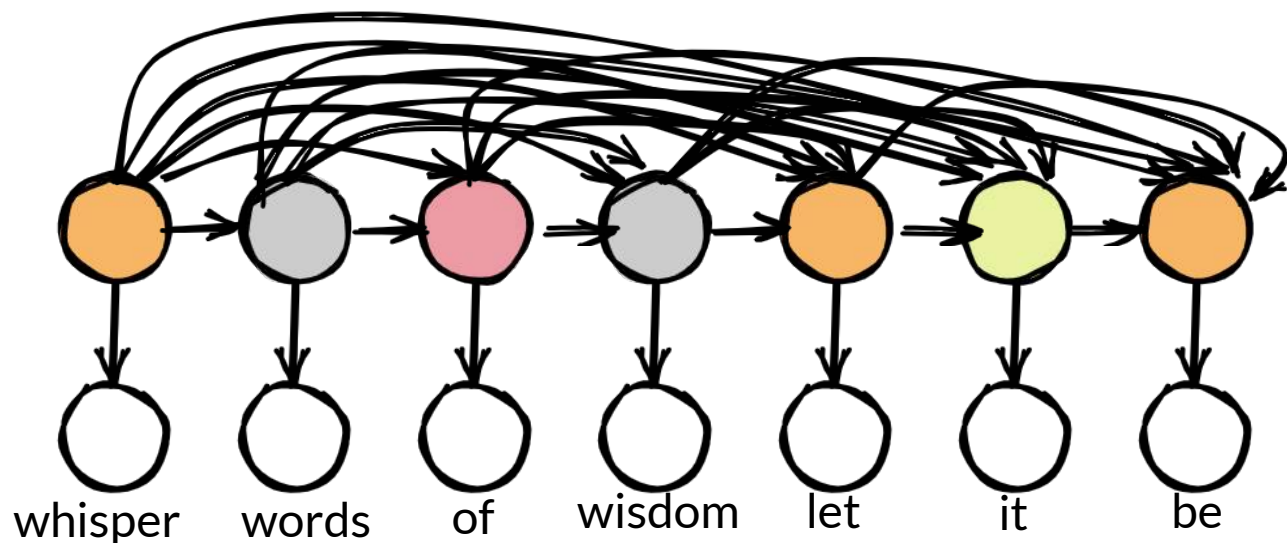$\times P(\text{words | noun}) \times P(\text{noun}) \times$
$\times \dots$
$= (0.7 \times 0.35) \times (0.7 \times 0.4) \times$
$\times (0.7 \times 0.15) \times \dots$

Emission model

# Better model

The previous probabilistic model is too simple. Context (adjacent words and labels) is essential. We'll add dependencies between labels (<u>not between words</u>)



whisper    words    of    wisdom    let    it    be

$$P(y, x) = P(x \mid y) \times P(y)$$
$$P(y) = P(y_1) \times P(y_2 \mid y_1) \times$$
$$\times P(y_3 \mid y_1, y_2) \times \ldots \times$$
$$\times P(y_N \mid y_1, y_2, \ldots, y_{N-1}) =$$

$$= P(y_1) \times \prod_{N=2} P(y_i \mid y_1, \ldots, y_{i-1})$$

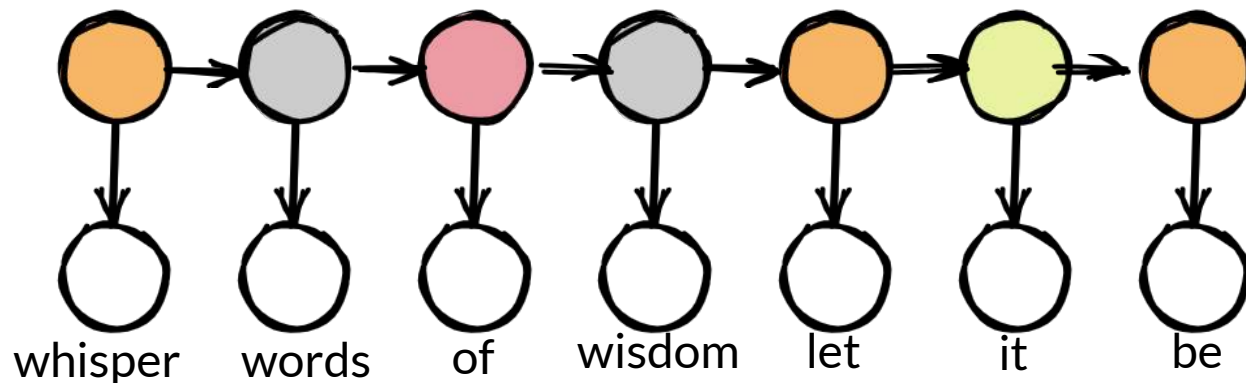thus each $y_i$ depends on all previous i−1 states

11

# Markov assumption

Hard to believe the $y_i$ element depends on all, including the first one.
Markov assumption allows us to consider $y_i$ being dependent on **only** the
last element ($y_{i-1}$)

$$P(y, x) = P(x \mid y) \times P(y) \, P(y)$$
$$= P(y_1) \times P(y_2 \mid y_1) \times$$
$$\times P(y_3 \mid y_2) \times \ldots \times P(y_N \mid y_{N-1}) =$$

$$= P(y_1) \times \prod_{N=2} P(y_i \mid y_{i-1})$$



whisper    words    of    wisdom    let    it    be

12

# Markov chain



state transition probability $P(y_i \mid y_{i-1})$

state $y_i$

whisper    words    of    wisdom    let    it    be

13

# Markov model

whisper words of wisdom let it be

| i | whispers | words | of | wisdom | let | it | be | v | n | p | p |
|---|---|---|---|---|---|---|---|---|---|---|---|
| verb (0.35) | **0.7** | 0.2 | 0.1 | 0.05 | **0.6** | 0.0 | **0.9** | 0.1 | 0.4 | 0.2 | 0.3 |
| noun (0.4) | 0.3 | **0.7** | 0.1 | **0.85** | 0.3 | 0.15 | 0.0 | 0.8 | 0.1 | 0.1 | 0.0 |
| prep (0.15) | 0.0 | 0.0 | **0.7** | 0.0 | 0.05 | 0.1 | 0.1 | 0.2 | 0.3 | 0.2 | 0.3 |
| pronoun (0.1) | 0.0 | 0.1 | 0.1 | 0.1 | 0.05 | **0.65** | 0.0 | 0.2 | 0.8 | 0.0 | 0.0 |

i-1

whisper   words   of   wisdom   let   it   be

$P(y, x) =$
$= P(\text{whisper} \mid \text{verb}) \times P(\text{verb}) \times$
$\times P(\text{words} \mid \text{noun}) \times P(\text{noun} \mid \text{verb}) \times$
$\times P(\text{of} \mid \text{prep}) \times P(\text{prep} \mid \text{noun}) \times$
$\times \ldots =$
$= (0.7 \times 0.35) \times (0.7 \times 0.8) \times$
$\times (0.7 \times 0.3) \times \ldots$

State transition matrix.
If constant => homogeneous
Markov model

14

# Hidden Markov Model

Hidden state = unobserved

observations

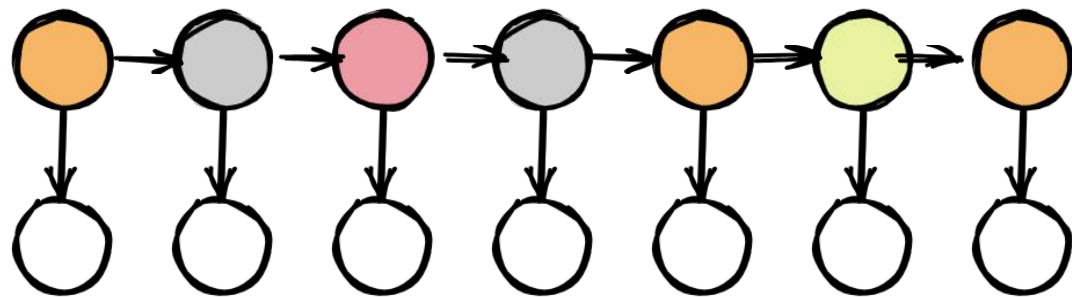whisper    words    of    wisdom    let    it    be

For a hidden Markov Model:
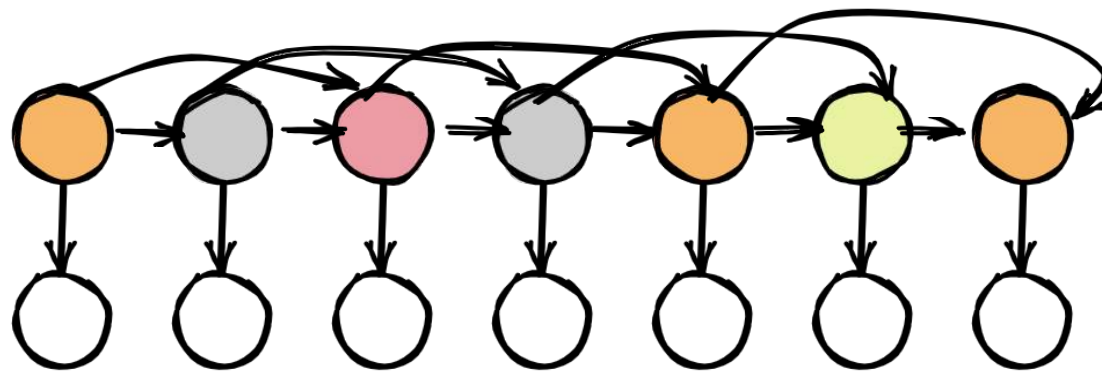
$$P(y, x) = P(x \mid y) \times P(y)$$

$P(y)$ is the probability of the hidden sequence. For a Markov chain $y_i$ depends only on previous state.

$$P(y) = i(y_1) \times \prod_{N=2} s(y_i, y_{i-1})$$

$P(x \mid y)$ is the emission model of the HMM  => $P(y, x)$ = Markov chain × emission model

15

# Higher-order HMMs

1st order HMM
bigram HMM

2nd order HMM
trigram HMM

# Inference problems for HMMs

Given an observation sequence $x$ and an HMM model $\lambda$, how do we efficiently compute $P(x|\lambda)$, i.e., the probability of the observation sequence given the model

Given an observation sequence $x$ and an HMM model $\lambda$, how do we choose a corresponding state sequence $y$ which is optimal in some sense, i.e., best explains the observations

Given an observation sequence $x$, how do we adjust (learn) the model parameters $\lambda$, to maximise $P(x|\lambda)$

Evaluation

forward algorithm

Decoding
(Recognition)

Viterbi algorithm

Training

Baum-Welch algorithm

17