# 10-301/601: Introduction to Machine Learning Lecture 23: Hidden Markov Models

Henry Chai

7/19/23

# Front Matter

- Announcements

  - PA5 released 7/13, due 7/20 (tomorrow) at 11:59 PM

  - PA6 released 7/20 (tomorrow), due 7/27 at 11:59 PM

- Recommended Readings

  - Murphy, Chapters 17.1 - 17.5

# Recall: Hidden Markov Models

- Two types of variables: observations (observed) and states (hidden or latent)
  - Set of states usually pre-specified via domain expertise/prior knowledge: $\{s_1, \dots, s_M\}$
  - Emission model:
    - Current observation is conditionally independent of all other variables given the current state: $P(X_t|Y_t)$
  - Transition model (Markov assumption):
    - Current state is conditionally independent of all earlier states given the previous state:
      $$P(Y_t|Y_{t-1}, \dots, Y_0) = P(Y_t|Y_{t-1})$$

# Hidden Markov Models: Outline

- How can we learn the conditional probabilities used by a hidden Markov model?

- What inference questions can we answer with a hidden Markov model?
    1. Computing the distribution of a single state (or a sequence of states) given a sequence of observations
    2. Finding the most-probable sequence of states given a sequence of observations
    3. Computing the probability of a sequence of observations

$$P(Y) = P(Y_1 \cap Y_2 \ldots Y_T)$$

# 3 Inference Questions for HMMs

1. Marginal Computation: $P\left(Y_t = s_j \mid \boldsymbol{x}^{(n)}\right)$ (or $P\left(Y \mid \boldsymbol{x}^{(n)}\right)$)

$$P\left(Y \mid \boldsymbol{x}^{(n)}\right) = \frac{P\left(\boldsymbol{x}^{(n)} \mid Y\right) P(Y)}{P\left(\boldsymbol{x}^{(n)}\right)} = \frac{\prod_{t=1}^{T} P\left(\boldsymbol{x}_t^{(n)} \mid Y_t\right) P(Y_t \mid Y_{t-1})}{P\left(\boldsymbol{x}^{(n)}\right)}$$

2. Decoding: $\hat{Y} = \underset{Y}{\operatorname{argmax}} \; P\left(Y \mid \boldsymbol{x}^{(n)}\right)$

3. Evaluation: $P\left(\boldsymbol{x}^{(n)}\right)$

$$P\left(\boldsymbol{x}^{(n)}\right) = \sum_{\mathcal{Y} \in \{\text{all possible sequences}\}} P\left(\boldsymbol{x}^{(n)} \mid \mathcal{Y}\right) P(\mathcal{Y})$$

of states

Sum rule of probability $= \sum \prod_{t=1}^{T} P\left(x_t^{(n)} \mid y_t\right) P(y_t \mid y_{t-1})$

# The Brute Force Algorithm

- Inputs: query $P\left(\boldsymbol{x}^{(n)}\right)$, emission matrix $A$, transition matrix $B$

- Initialize $p = 0$

- For $\mathcal{Y} \in \{\text{all possible sequences}\} \rightarrow M^T$

  *of states*

  - Compute the joint probability

$$P\left(\boldsymbol{x}^{(n)}, \mathcal{Y}\right) = P\left(\boldsymbol{x}^{(n)} | \mathcal{Y}\right) P(\mathcal{Y}) = \prod_{t=1}^{T} P\left(\boldsymbol{x}_t^{(n)} \Big| \mathcal{Y}_t\right) P(\mathcal{Y}_t | \mathcal{Y}_{t-1})$$

  - $p \mathrel{+}= P\left(\boldsymbol{x}^{(n)}, \mathcal{Y}\right)$

- Return $p = P\left(\boldsymbol{x}^{(n)}\right)$

# Lecture 23 Polls

**0 done**

↻ **0 underway**

# Given $C$ possible observations and $M$ possible states plus special START/END states, how many possible sequences of length $T$ (not counting the START and END states) are there?
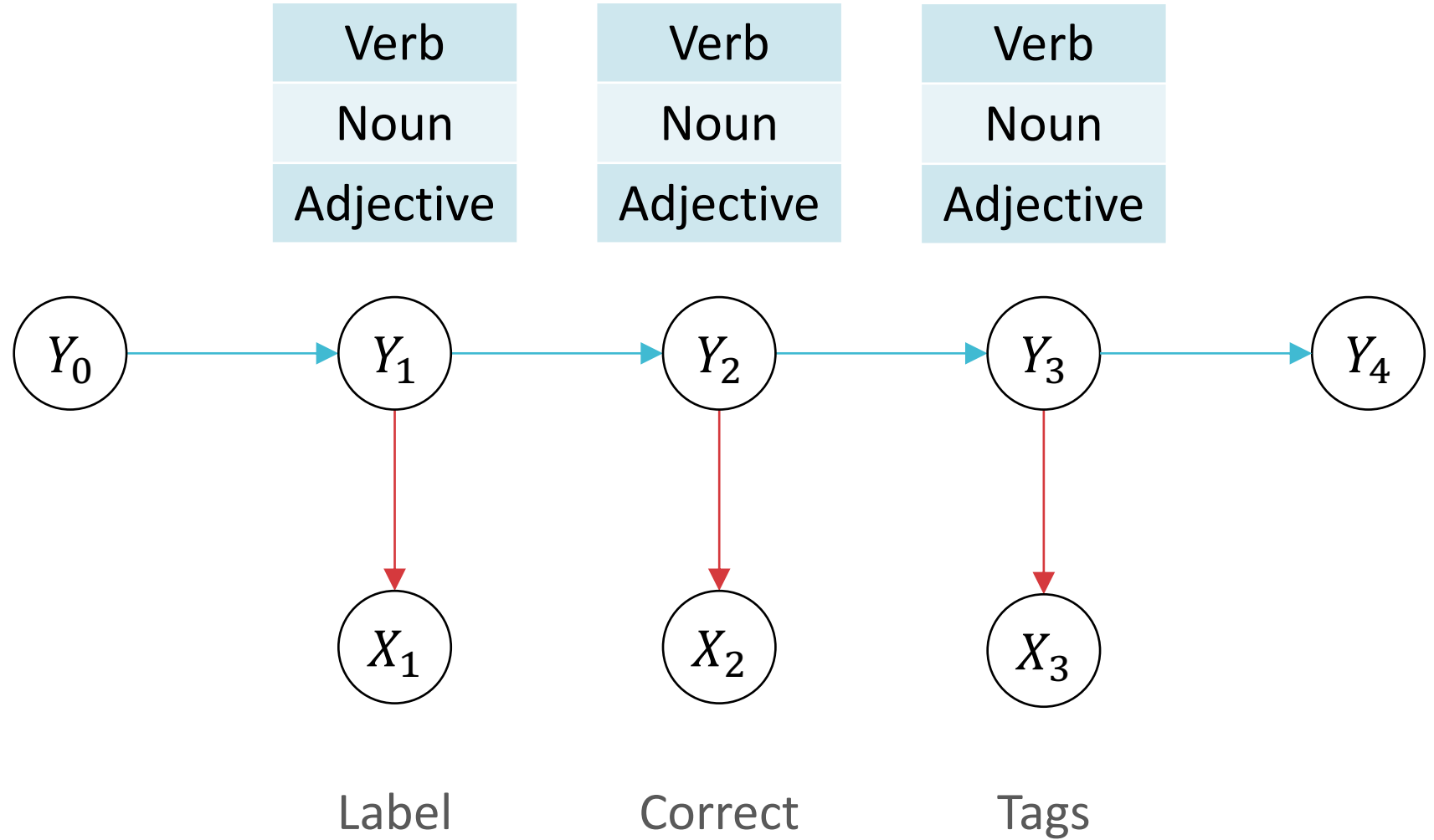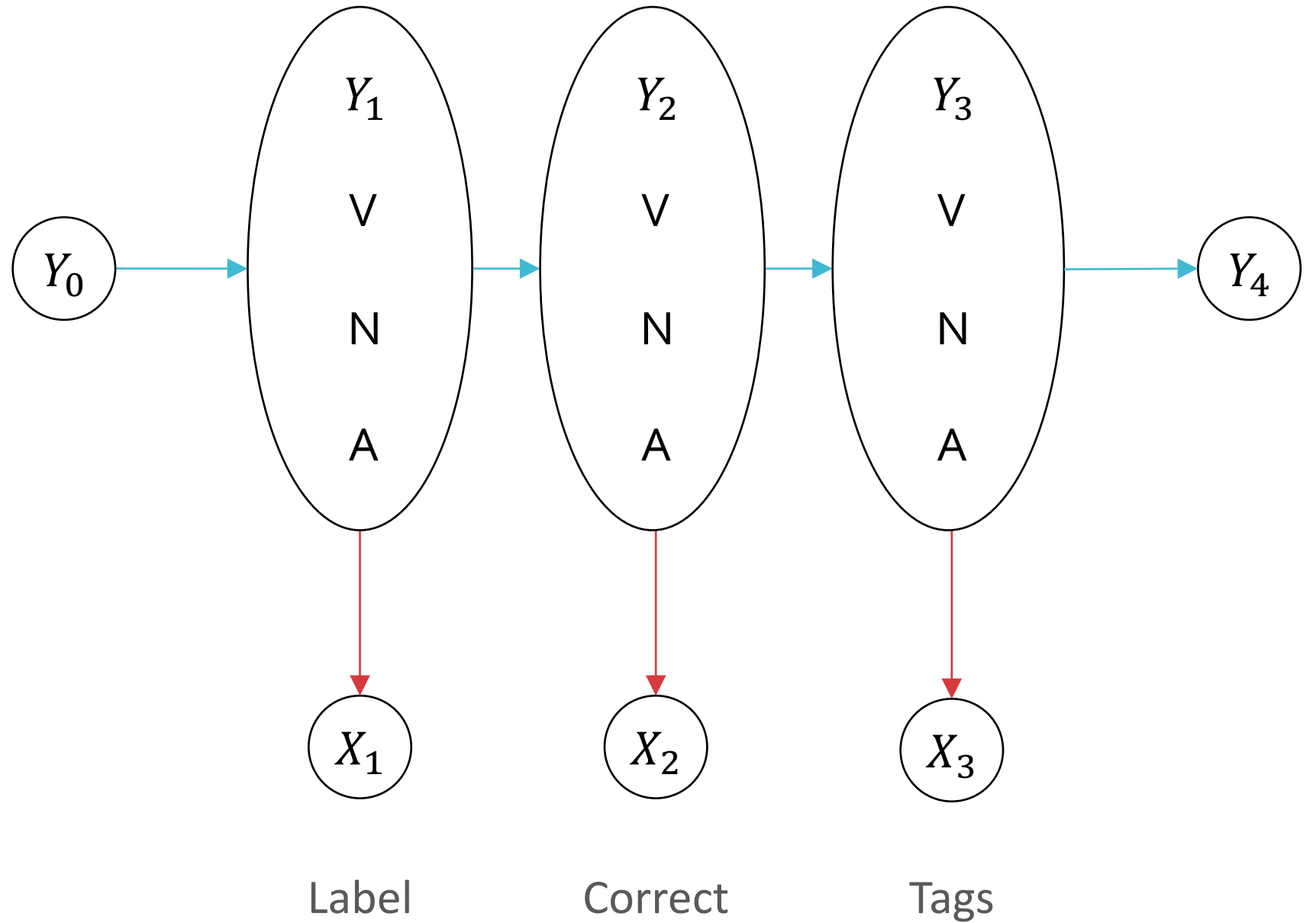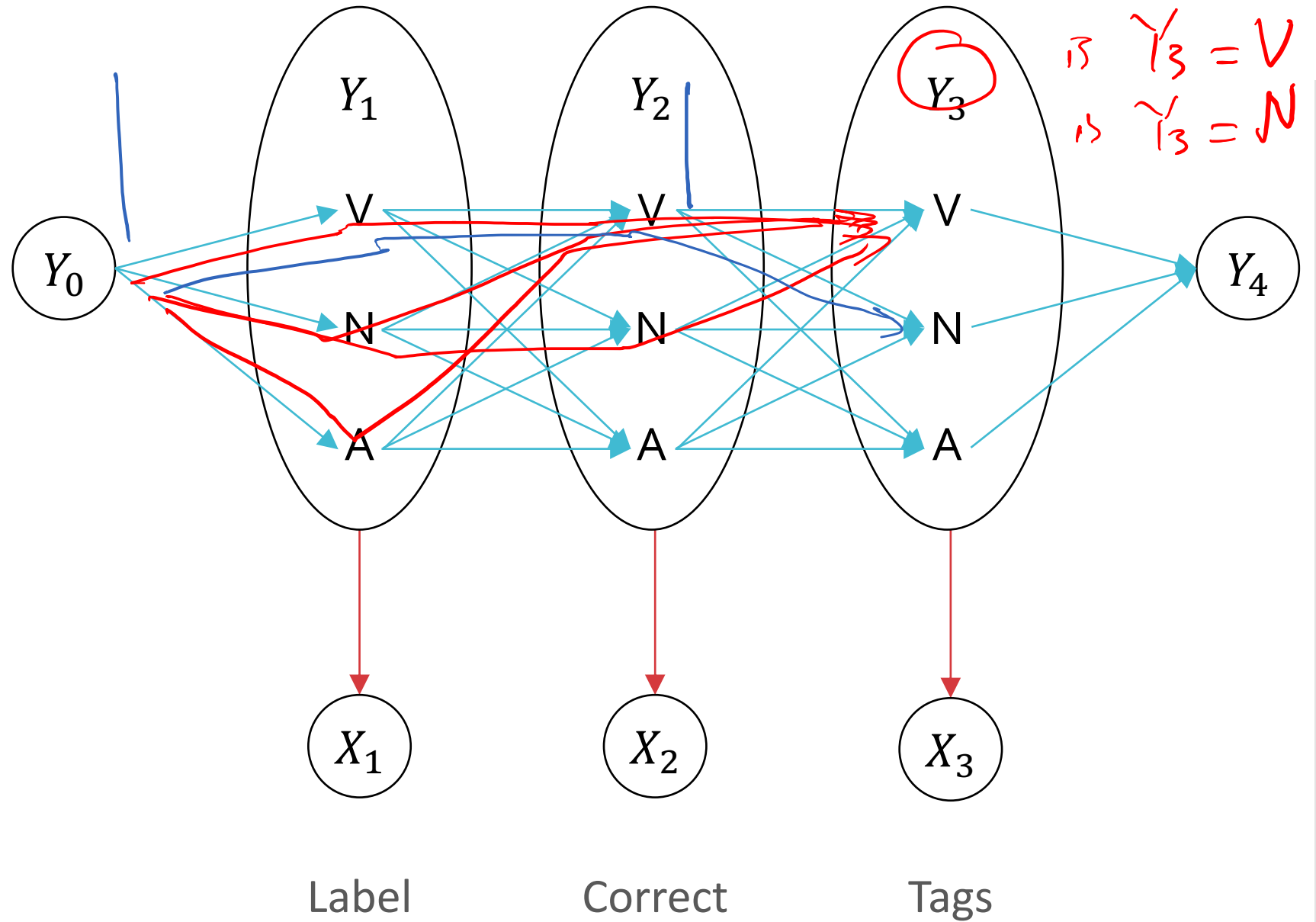
$TC$

$TM$

$T^M$

$M^T$

# Inference with HMMs: PoS Tagging Example

| Verb | | Verb | | Verb | |
|------|--|------|--|------|--|
| Noun | | Noun | | Noun | |
| Adjective | | Adjective | | Adjective | |

$Y_0 \rightarrow Y_1 \rightarrow Y_2 \rightarrow Y_3 \rightarrow Y_4$

$Y_1 \rightarrow X_1$

$Y_2 \rightarrow X_2$

$Y_3 \rightarrow X_3$

Label    Correct    Tags

Inference with HMMs: PoS Tagging Example

# Inference with HMMs: PoS Tagging Example



Label          Correct          Tags

# 3 Inference Questions for HMMs

1. Marginal Computation: $P\left(Y_t = s_j \,\middle|\, \boldsymbol{x}^{(n)}\right)$ (or $P\left(Y \,\middle|\, \boldsymbol{x}^{(n)}\right)$)

$$P\left(Y_t = s_j \,\middle|\, \boldsymbol{x}^{(n)}\right) = \frac{P\left(Y_t = s_j, \boldsymbol{x}^{(n)}\right)}{P\left(\boldsymbol{x}^{(n)}\right)}$$

2. Decoding: $\hat{Y} = \underset{Y}{\operatorname{argmax}} \; P\left(Y \,\middle|\, \boldsymbol{x}^{(n)}\right)$

3. Evaluation: $P\left(\boldsymbol{x}^{(n)}\right)$

$$P\left(\boldsymbol{x}^{(n)}\right) = \sum_{m=1}^{M} P\left(Y_t = s_m, \boldsymbol{x}^{(n)}\right)$$

# Recursive Marginals

$$P\left(Y_t = s_j, \boldsymbol{x}_1^{(n)}, \ldots, \boldsymbol{x}_T^{(n)}\right)$$

$$= P\left(x_{t+1}^{(n)}, \ldots, x_T^{(n)} \mid Y_t = s_j, x_1^{(n)}, \ldots, x_t^{(n)}\right) P\left(Y_t = s_j, x_1^{(n)}, \ldots, x_t^{(n)}\right)$$

$$\downarrow \text{ by conditional independences of HMMs}$$

$$= P\left(x_{t+1}^{(n)}, \ldots, x_T^{(n)} \mid Y_t = s_j\right) P\left(Y_t = s_j, x_1^{(n)}, \ldots, x_t^{(n)}\right)$$

$$:= \beta_t(j)\, \alpha_t(j)$$

## Forward Algorithm

$$\alpha_t(j) := P\left(Y_t = s_j, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_t^{(n)}\right)$$

$$= \sum_{m=1}^{M} P\left(Y_t = s_j, Y_{t-1} = s_m, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_t^{(n)}\right)$$

$$= \sum_{m=1}^{M} P\left(\boldsymbol{x}_t^{(n)} \mid Y_t = s_j, Y_{t-1} = s_m, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_{t-1}^{(n)}\right) *$$

$$P\left(Y_t = s_j, Y_{t-1} = s_m, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_{t-1}^{(n)}\right)$$

# Forward Algorithm

$$\alpha_t(j) := P\left(Y_t = s_j, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_t^{(n)}\right)$$

Base case:
$$\alpha_0(\text{START}) = 1$$
$$\alpha_0(s_m) = 0$$

$$= \sum_{m=1}^{M} P\left(Y_{t-1} = s_m, Y_t = s_j, x_1^{(n)}, \dots, x_t^{(n)}\right)$$

$$= \sum_{m=1}^{M} P\left(x_t^{(n)} \mid Y_{t-1} = s_m, Y_t = s_j, x_1^{(n)}, \dots, x_{t-1}^{(n)}\right) \cdot$$
$$\left( P\left(Y_{t-1} = s_m, Y_t = s_j, x_1^{(n)}, \dots, x_{t-1}^{(n)}\right) \right.$$

$$= \sum_{m=1}^{M} P\left(x_t^{(n)} \mid Y_t = s_j\right) P\left(Y_t = s_j \mid Y_{t-1} = s_m, x_1^{(n)}, \dots, x_{t-1}^{(n)}\right)$$
$$P\left(Y_{t-1} = s_m, x_1^{(n)}, \dots, x_{t-1}^{(n)}\right)$$

$$= \sum_{m=1}^{M} P\left(x_t^{(n)} \mid Y_t = s_j\right) \frac{P\left(Y_t = s_j \mid Y_{t-1} = s_m\right)}{P\left(Y_{t-1} = s_m, x_1^{(n)}, \dots, x_{t-1}^{(n)}\right) = \alpha_{t-1}^{(m)}}$$

$$= \sum_{m=1}^{M} P\left(x_t^{(n)} \mid Y_t = s_j\right) P\left(Y_t = s_j \mid Y_{t-1} = s_m\right) \alpha_{t-1}(m)$$

# Backward Algorithm

$$\beta_t(j) := P\left(\boldsymbol{x}_{t+1}^{(n)}, \ldots, \boldsymbol{x}_T^{(n)} \middle| Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(\boldsymbol{x}_{t+1}^{(n)}, \ldots, \boldsymbol{x}_T^{(n)}, Y_{t+1} = s_m \middle| Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(\boldsymbol{x}_{t+2}^{(n)}, \ldots, \boldsymbol{x}_T^{(n)} \middle| Y_t = s_j, \boldsymbol{x}_{t+1}^{(n)}, Y_{t+1} = s_m\right) *$$

$$P\left(\boldsymbol{x}_{t+1}^{(n)}, Y_{t+1} = s_m \middle| Y_t = s_j\right)$$

# Backward Algorithm

$$\beta_t(j) := P\left(x_{t+1}^{(n)}, \ldots, x_T^{(n)} \mid Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(x_{t+1}^{(n)}, \ldots, x_T^{(n)}, Y_{t+1} = s_m \mid Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(x_{t+2}^{(n)}, \ldots, x_T^{(n)} \mid x_{t+1}^{(n)}, Y_{t+1} = s_m, Y_t = s_j\right) \cdot$$

$$P\left(x_{t+1}^{(n)}, Y_{t+1} = s_m \mid Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(x_{t+2}^{(n)}, \ldots, x_T^{(n)} \mid Y_{t+1} = s_m\right) \cdot$$

$$P\left(x_{t+1}^{(n)} \mid Y_{t+1} = s_m, Y_t = s_j\right) P\left(Y_{t+1} = s_m \mid Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} P\left(x_{t+2}^{(n)}, \ldots, x_T^{(n)} \mid Y_{t+1} = s_m\right) P\left(x_{t+1}^{(n)} \mid Y_{t+1} = s_m\right)$$

$$P\left(Y_{t+1} = s_m \mid Y_t = s_j\right)$$

$$= \sum_{m=1}^{M} \beta_{t+1}(m) P\left(Y_{t+1} = s_m \mid Y_t = s_j\right) P\left(x_{t+1}^{(n)} \mid Y_{t+1} = s_m\right)$$

# The Forward-Backward Algorithm

- Inputs: query $P\left(Y_t = s_j \mid \boldsymbol{x}^{(n)}\right)$, emission matrix $A$, transition matrix $B$

- Initialize $\alpha_0(\text{START}) = 1$ and $\beta_{T+1}(\text{END}) = 1$

- For $\tau = 1, \ldots, T$
  - For $m = 1, \ldots, M$

$$\alpha_\tau(m) = P\left(\boldsymbol{x}_\tau^{(n)} \mid Y_\tau = s_m\right) \sum_{k=1}^{M} P(Y_\tau = s_m \mid Y_{\tau-1} = s_k)\alpha_{\tau-1}(k)$$

$$TMM = O(TM^2)$$

- For $\tau = T, \ldots, 1$
  - For $m = 1, \ldots, M$

$$\beta_\tau(m) = \sum_{k=1}^{M} \beta_{\tau+1}(k) P\left(\boldsymbol{x}_{\tau+1}^{(n)} \mid Y_{\tau+1} = s_k\right) P(Y_{\tau+1} = s_k \mid Y_\tau = s_m)$$

- Return $P\left(Y_t = s_j \mid \boldsymbol{x}^{(n)}\right) = \dfrac{P\left(Y_t = s_j, \boldsymbol{x}^{(n)}\right)}{P\left(\boldsymbol{x}^{(n)}\right)} = \dfrac{\beta_t(j)\alpha_t(j)}{\sum_{m=1}^{M} \beta_t(m)\alpha_t(m)}$

$$P\left(\boldsymbol{x}^{(n)}\right) = \sum_{m=1}^{M} P\left(Y_t = s_m, \boldsymbol{x}^{(n)}\right)$$

Given $C$ possible observations and $M$ possible states plus special START/END states, what is the runtime of the forward-backward algorithm on sequences of length $T$?
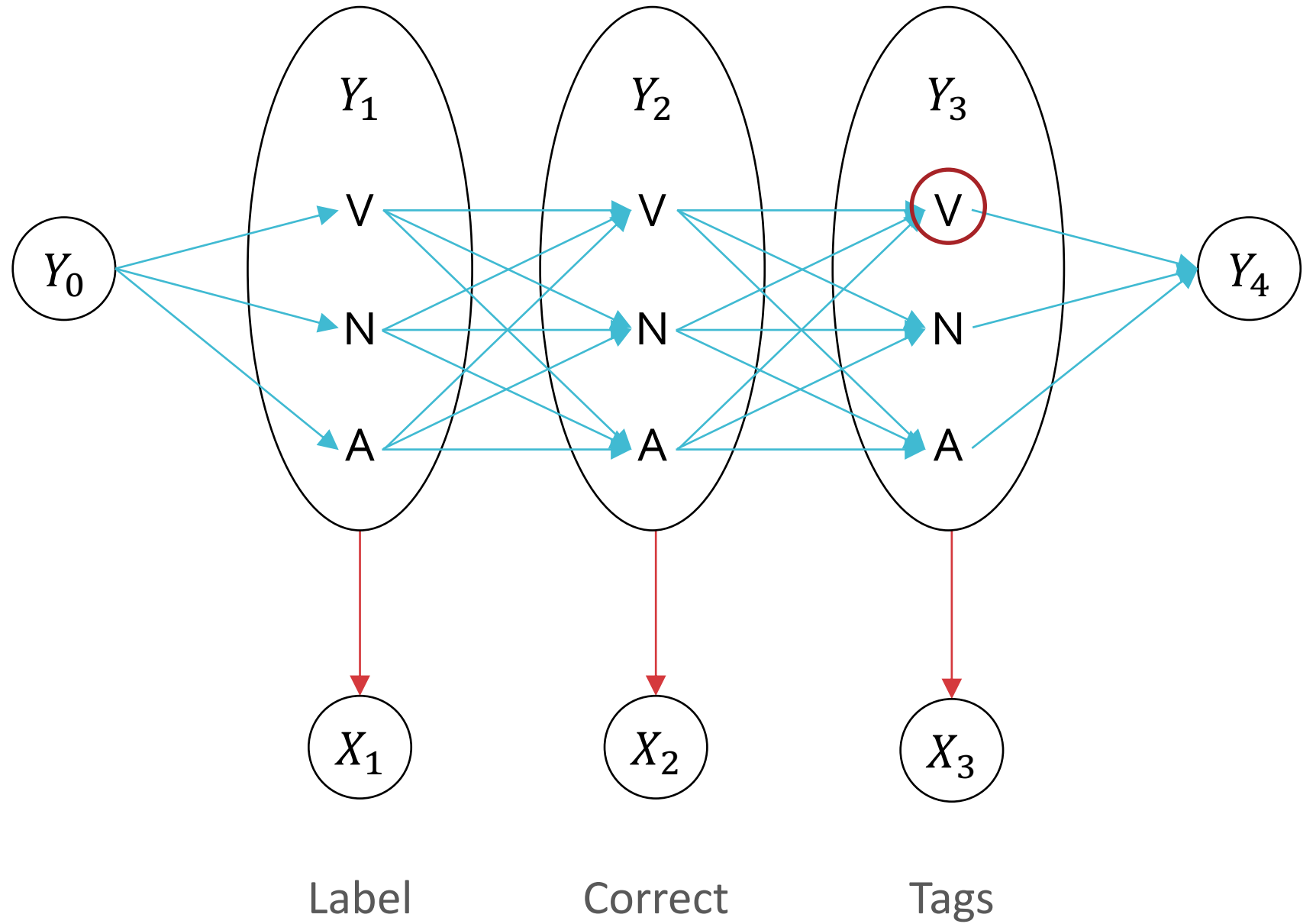
$O(TM)$

$O(T^2 M)$

$O(TM^2)$

$O(T^2 M^2)$

# Most Probable State Sequence

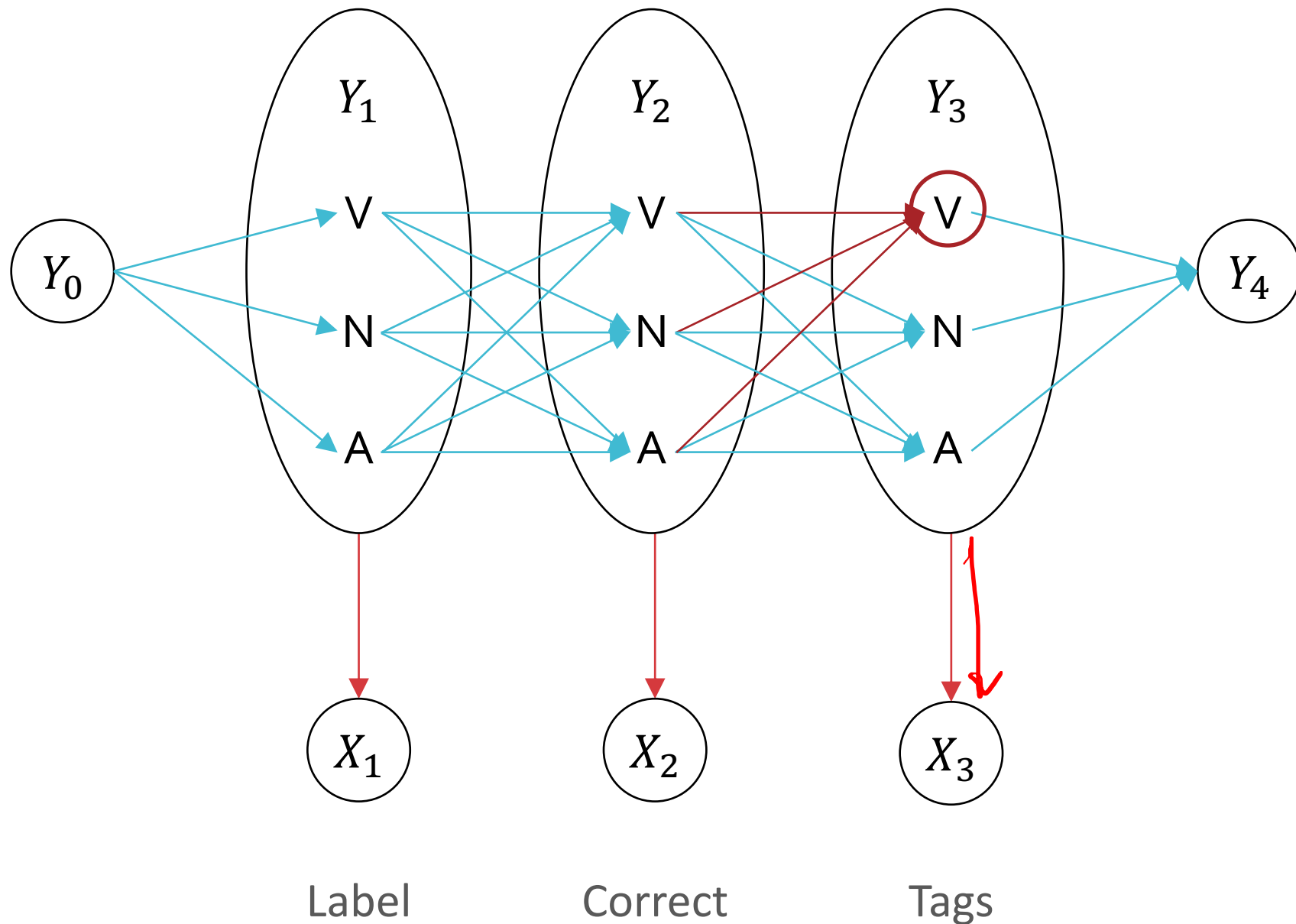Decoding: $\hat{Y} = \underset{Y}{\text{argmax}} \; P(Y|\boldsymbol{x}^{(n)})$

$$\omega_t(j) := \underset{\mathcal{Y} \in \{\text{all possible sequences of } t-1 \text{ states}\}}{\max} P\left(\mathcal{Y}, Y_t = s_j, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_t^{(n)}\right)$$

$=$ the probability of the most probable sequence of $t$ states that ends in $s_j$, conditioned on the first $t$ observations
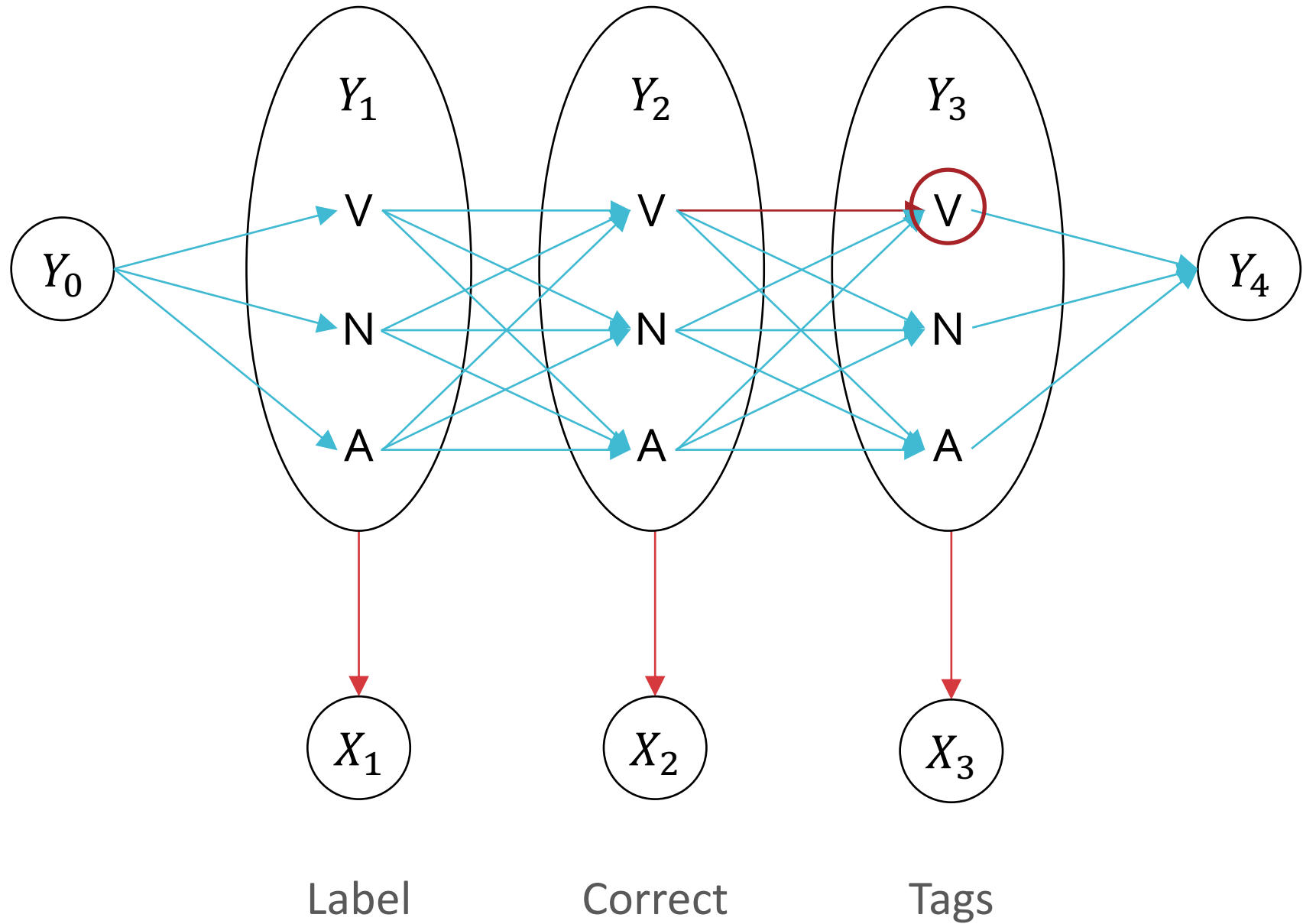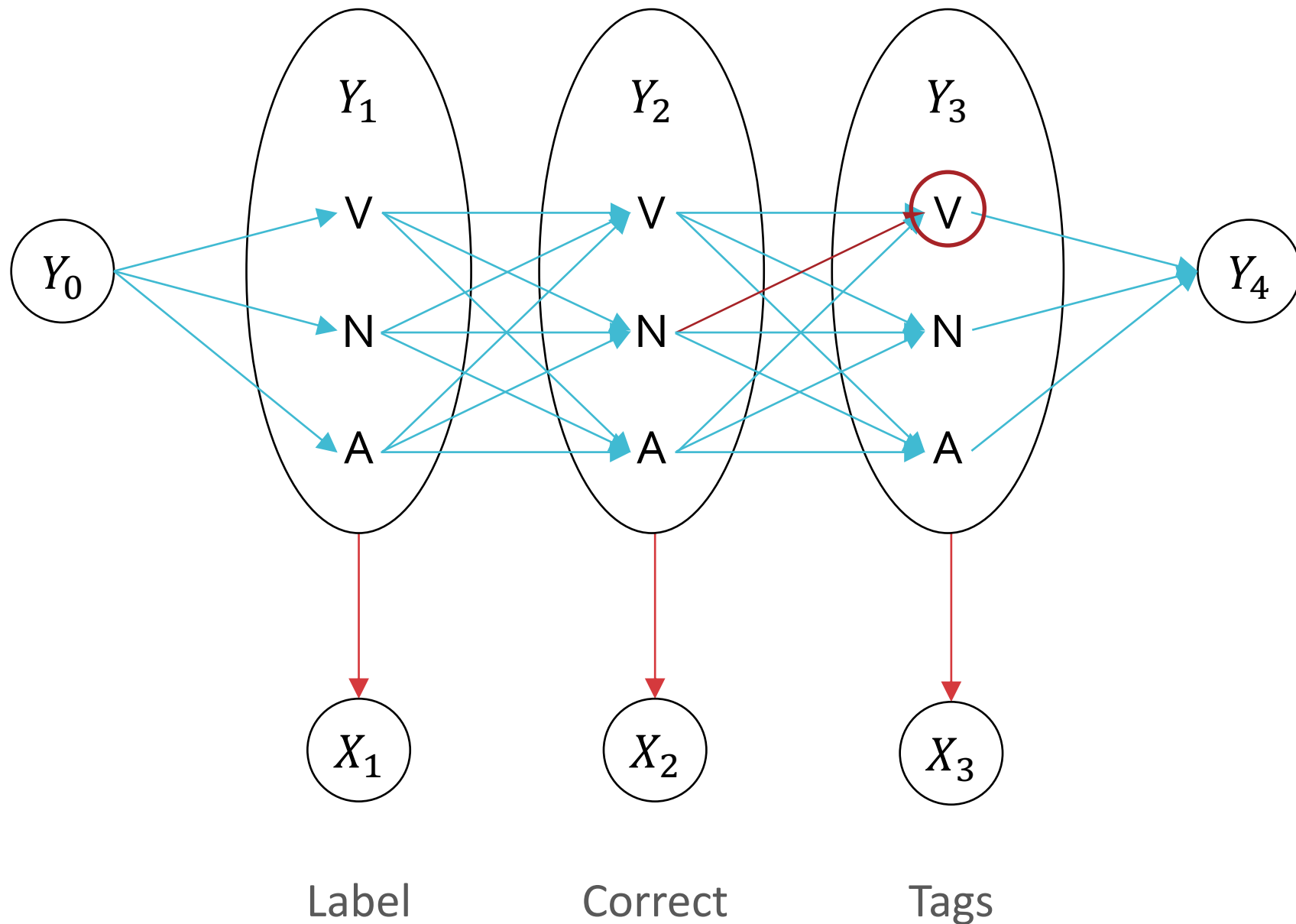
$\omega_3(V)$

Label        Correct        Tags

$$\omega_3(V)$$
$$= \max\{$$
$$\omega_2(V) P(Y_3 = V \mid Y_2 = V)$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\omega_2(N) P(Y_3 = V \mid Y_2 = N)$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\omega_2(A) P(Y_3 = V \mid Y_2 = A)$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\}$$



Label      Correct      Tags

$$\omega_3(V)$$
$$= \max\{$$
$$\omega_2({\color{red}V})P(Y_3 = V | Y_2 = {\color{red}V})$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\omega_2(N)P(Y_3 = V | Y_2 = N)$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\omega_2(A)P(Y_3 = V | Y_2 = A)$$
$$P\left(x_3^{(n)} \middle| Y_3 = V\right),$$
$$\}$$

Label          Correct          Tags

$$\omega_3(V)$$
$$= \max\{$$
$$\omega_2(V)P(Y_3 = V | Y_2 = V)$$
$$P\left(\boldsymbol{x}_3^{(n)} \Big| Y_3 = V\right),$$
$$\omega_2(N)P(Y_3 = V | Y_2 = N)$$
$$P\left(\boldsymbol{x}_3^{(n)} \Big| Y_3 = V\right),$$
$$\omega_2(A)P(Y_3 = V | Y_2 = A)$$
$$P\left(\boldsymbol{x}_3^{(n)} \Big| Y_3 = V\right),$$
$$\}$$

Label          Correct          Tags

$\omega_3(V)$

$= \max\{$

$\omega_2(V)P(Y_3 = V | Y_2 = V)$

$P\left(x_3^{(n)} \middle| Y_3 = V\right),$

$\omega_2(N)P(Y_3 = V | Y_2 = N)$

$P\left(x_3^{(n)} \middle| Y_3 = V\right),$

$\omega_2(A)P(Y_3 = V | Y_2 = A)$

$P\left(x_3^{(n)} \middle| Y_3 = V\right),$

$\}$

# Most Probable State Sequence

Decoding: $\hat{Y} = \underset{Y}{\operatorname{argmax}} \, P\left(Y \middle| \boldsymbol{x}^{(n)}\right)$

$$\omega_t(j) := \underset{\mathcal{Y} \,\in\, \{\text{all possible sequences of } t-1 \text{ states}\}}{\max} P\left(\mathcal{Y}, Y_t = s_j, \boldsymbol{x}_1^{(n)}, \dots, \boldsymbol{x}_t^{(n)}\right)$$

= the probability of the most probable sequence of $t$ states that ends in $s_j$, conditioned on the first $t$ observations

$$= \underset{m \,\in\, \{1,\dots,M\}}{\max} \omega_{t-1}(m) \, P\left(Y_t = s_j \middle| Y_{t-1} = s_m\right) P\left(\boldsymbol{x}_t^{(n)} \middle| Y_t = s_j\right)$$

# The Viterbi Algorithm

- Inputs: observations $\boldsymbol{x}^{(n)}$, emission matrix $A$, transition matrix $B$

- Initialize $\omega_0(\text{START}) = 1$

- For $\tau = 1, \dots, T + 1$
    - For $m = 1, \dots, M$

$$\rightarrow \omega_\tau(m) = \max_{k \in \{1,\dots,M\}} P\left(\boldsymbol{x}_\tau^{(n)} \mid Y_\tau = s_m\right) P(Y_\tau = s_m \mid Y_{\tau-1} = s_k)\omega_{\tau-1}(k)$$

$$\rho_\tau(m) = \operatorname*{argmax}_{k \in \{1,\dots,M\}} P\left(\boldsymbol{x}_\tau^{(n)} \mid Y_\tau = s_m\right) P(Y_\tau = s_m \mid Y_{\tau-1} = s_k)\omega_{\tau-1}(k)$$

- Return the most probable assignment given $\boldsymbol{x}^{(n)}$:
    - $\hat{Y}_T = \rho_{T+1}(\text{END})$
    - For $\tau = T - 1, \dots, 1$
        - $\hat{Y}_\tau = \rho_{\tau+1}\left(\hat{Y}_{\tau+1}\right)$

## 3̶ 4 Inference Questions for HMMs

1. Marginal Computation: $P\left(Y_t = s_j \middle| \boldsymbol{x}^{(n)}\right)$ (or $P\left(Y \middle| \boldsymbol{x}^{(n)}\right)$)

$$P\left(Y \middle| \boldsymbol{x}^{(n)}\right) = \frac{P\left(\boldsymbol{x}^{(n)} \middle| Y\right)P(Y)}{P\left(\boldsymbol{x}^{(n)}\right)} = \frac{\prod_{t=1}^{T} P\left(\boldsymbol{x}_t^{(n)} \middle| Y_t\right) P(Y_t | Y_{t-1})}{P\left(\boldsymbol{x}^{(n)}\right)}$$

2. <u>Viterbi</u> Decoding: $\hat{Y} = \underset{Y}{\text{argmax}} \; P\left(Y \middle| \boldsymbol{x}^{(n)}\right)$

3. Evaluation: $P\left(\boldsymbol{x}^{(n)}\right)$

$$P\left(\boldsymbol{x}^{(n)}\right) = \sum_{\mathcal{Y} \in \{\text{all possible sequences}\}} P\left(\boldsymbol{x}^{(n)} \middle| \mathcal{Y}\right) P(\mathcal{Y})$$

4. Minimum Bayes Risk (MBR) Decoding:

$$\hat{Y} = \underset{Y}{\text{argmin}} \; \mathbb{E}_{Y' \sim P_{A,B}\left(\cdot \middle| \boldsymbol{x}^{(n)}\right)}[\ell(Y, Y')]$$

# Minimum Bayes Risk Decoding

- The learned parameters $A$ and $B$ induce a probability distribution or belief over sequences of states $P_{A,B}\left(Y\middle|\boldsymbol{x}^{(n)}\right)$

- Given a loss function, $\ell(Y, Y')$, find the sequence of states that minimizes our expected loss *under our current belief*

$$\hat{Y} = \underset{Y}{\text{argmin}} \ \mathbb{E}_{Y' \sim P_{A,B}\left(\cdot\middle|\boldsymbol{x}^{(n)}\right)}[\ell(Y, Y')]$$

$$= \underset{Y}{\text{argmin}} \ \sum_{Y'} P_{A,B}\left(Y'\middle|\boldsymbol{x}^{(n)}\right) \ell(Y, Y')$$

# Minimum Bayes Risk Decoding: Example

- If $\ell(Y, Y')$ is the 0-1 loss

$$\ell(Y, Y') = 1 - \mathbb{1}(Y = Y')$$

$$\hat{Y} = \underset{Y}{\mathrm{argmin}} \sum_{Y'} P_{A,B}(Y'|\boldsymbol{x}^{(n)})\left(1 - \mathbb{1}(Y = Y')\right)$$

$$= \underset{Y}{\mathrm{argmin}} -\sum_{Y'} P_{A,B}(Y'|\boldsymbol{x}^{(n)})\,\mathbb{1}(Y = Y')$$

$$= \underset{Y}{\mathrm{argmax}}\; P_{A,B}(Y|\boldsymbol{x}^{(n)})$$

# Minimum Bayes Risk Decoding: Example

- If $\ell(Y, Y')$ is the Hamming loss

$$\ell(Y, Y') = \sum_{t=1}^{T} 1 - \mathbb{1}(Y_t = Y_t')$$

$$\hat{Y}_t = \underset{Y_t}{\operatorname{argmax}} \; P_{A,B}\left(Y_t \middle| \boldsymbol{x}^{(n)}\right)$$

- Computes the most likely state at each time step using the marginals from the forward-backward algorithm

# Key Takeaways

- Because of their well-behaved graphical structure, inference in HMMs is tractable via dynamic programming
  - Forward-backward algorithm for computing marginal distributions
  - Viterbi algorithm for computing most probable sequence of states