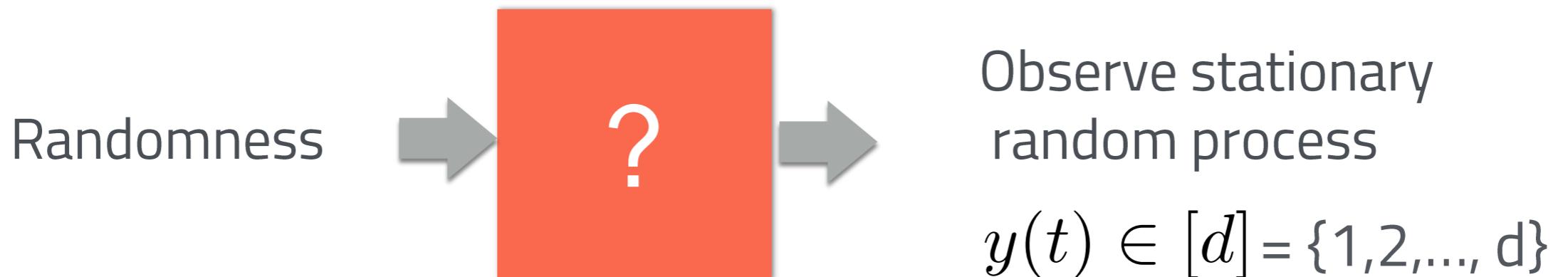


Minimal Realization Problems for Hidden Markov Models

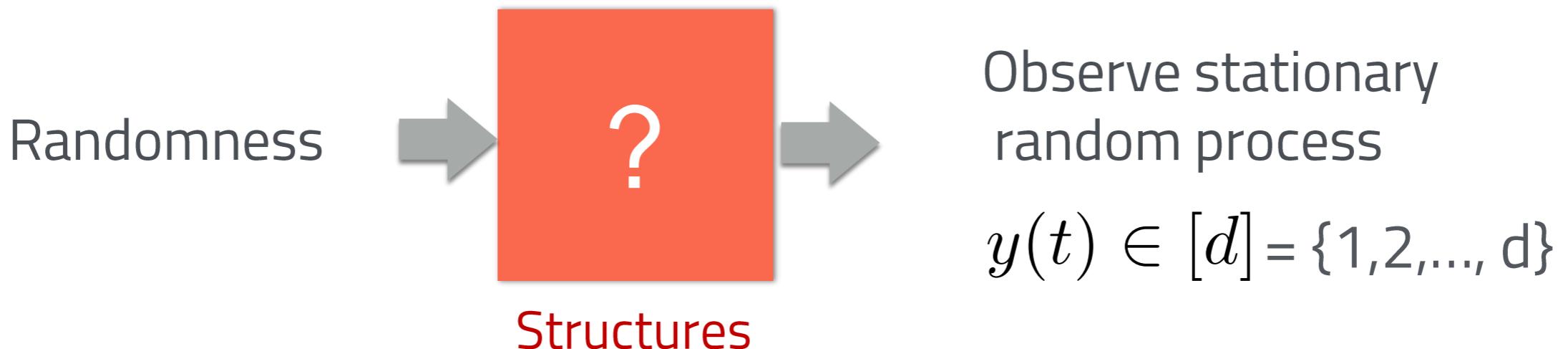
Motivation



- ♦ The random process is characterized by the joint distribution
(Probabilities of strings of any length N)

$$\mathcal{S}^{(\infty)} = \left\{ \mathbb{P}(y_1 = l_1, \dots, y_N = l_N) : \quad \forall \mathbf{l}_1^N \in [d]^N, \quad \forall N \in \mathbb{Z} \right\}$$

Motivation

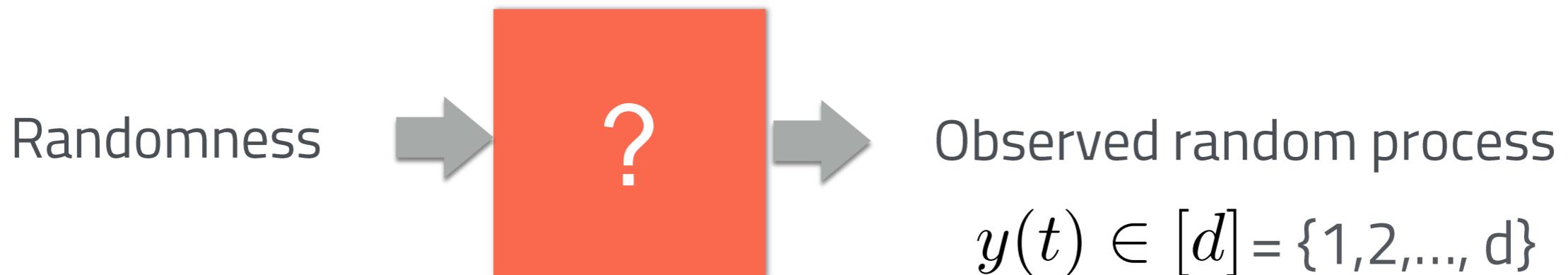


- ♦ The random process is characterized by the joint distribution
(Probabilities of strings of any length N)

$$\mathcal{S}^{(\infty)} = \left\{ \mathbb{P}(y_1 = l_1, \dots, y_N = l_N) : \quad \forall \mathbf{l}_1^N \in [d]^N, \quad \forall N \in \mathbb{Z} \right\}$$

- ♦ Can we find a **finite state system** description of the process?

Motivation



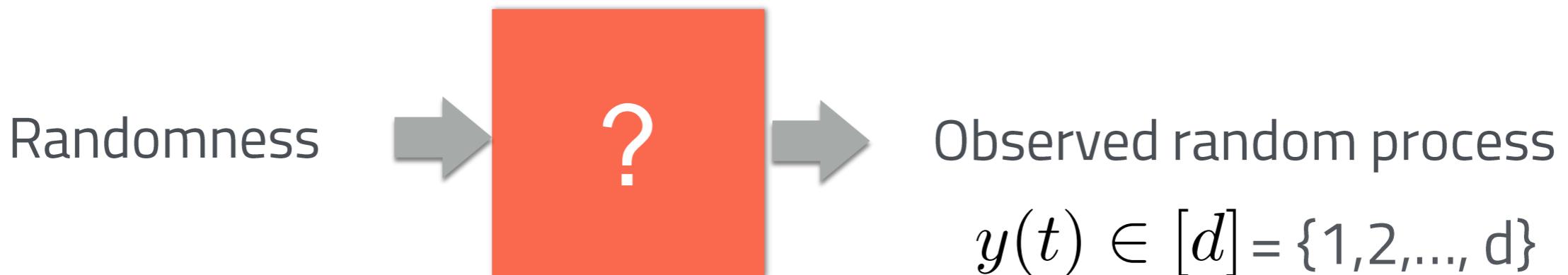
- ♦ Assume that the process can be described by a finite state system

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

$$w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v,$$
$$z_t = u' w_t, \quad \forall \mathbf{l}_1^N \in [d]^N, \forall N \in \mathbb{Z}$$

$$\begin{aligned} \mathbb{P}(y_1 = l_1, \dots, y_N = l_N) &= z_N(\mathbf{l}_1^N) \\ &= u' A^{(l_n)} \dots A^{(l_1)} A^{(l_{-1})} \dots A^{(l_{-n})} v \end{aligned}$$

Motivation



- Assume that the process can be described by a finite state system

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

$$w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v, \\ z_t = u' w_t,$$

State has no physical meaning

$$w_t \in \mathbb{R}^k$$

$$\forall \mathbf{l}_1^N \in [d]^N, \forall N \in \mathbb{Z}$$

$$\mathbb{P}(y_1 = l_1, \dots, y_N = l_N) = z_N(\mathbf{l}_1^N)$$

$$= u' A^{(l_n)} \dots A^{(l_1)} A^{(l_{-1})} \dots A^{(l_{-n})} v$$

Finite state system realization

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

- ♦ **Multi-linear system** $w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v,$
 $z_t = u' w_t,$

Finite state system realization

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

♦ Multi-linear system $w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v,$

$$z_t = u' w_t,$$

♦ Quasi-HMM realization

If $\mathbb{P}(y_1 = l_1, \dots, y_N = l_N) = z_N(\mathbf{l}_1^N)$ is a valid joint distribution

and $u' \left[\sum_{j \in [d]} A^{(j)} \right] = u', \quad \left[\sum_{j \in [d]} A^{(j)} \right] v = v.$ (stationarity)

Finite state system realization

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

♦ Multi-linear system $w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v,$

$$z_t = u' w_t,$$

♦ Quasi-HMM realization

If $\mathbb{P}(y_1 = l_1, \dots, y_N = l_N) = z_N(\mathbf{l}_1^N)$ is a valid joint distribution

and $u' \left[\sum_{j \in [d]} A^{(j)} \right] = u', \quad \left[\sum_{j \in [d]} A^{(j)} \right] v = v.$ (stationarity)

♦ HMM realization

If in addition $v = \mathbf{e}_k, \quad u \in \mathbb{R}_+^k, \quad \sum_i u_i = 1,$

$$A^{(j)} \in [0, 1]^{k \times k}, \forall j \in [d],$$

Finite state system realization

$$\theta = (d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

♦ Multi-linear system $w_{t+1} = A^{(l_t)} w_t, \quad w_0 = v,$

$$z_t = u' w_t,$$

♦ Quasi-HMM realization

If $\mathbb{P}(y_1 = l_1, \dots, y_N = l_N) = z_N(\mathbf{l}_1^N)$ is a valid joint distribution

and $u' \left[\sum_{j \in [d]} A^{(j)} \right] = u', \quad \left[\sum_{j \in [d]} A^{(j)} \right] v = v.$ (stationarity)

♦ HMM realization

If in addition $v = \mathbf{e}_k, \quad u \in \mathbb{R}_+^k, \quad \sum_i u_i = 1,$

$$A^{(j)} \in [0, 1]^{k \times k}, \forall j \in [d],$$

1-1 mapping to

Transition: $Q \in \mathbb{R}^{k \times k}$

Observation: $O \in \mathbb{R}^{d \times k}$

Realization problems

- ♦ **Input:** probabilities of length N strings of an HMM

$$\mathcal{S}^{(N)} = \{\mathbb{P}(y_1 = l_1, \dots, y_N = l_N) : \forall \mathbf{l}_1^N \in [d]^N\}$$

- ♦ **Output:**

order

$$(d, k, u \in \mathbb{R}^k, v \in \mathbb{R}^k, A^{(l)} \in \mathbb{R}^{k \times k} : l \in [d])$$

✓

Minimal order quasi-HMM realization

✓

Minimal order HMM realization

$$\mathcal{S}^{(\infty)} = \left\{ \mathbb{P}(\mathbf{y}_1^N = \mathbf{l}_1^N) : \forall \mathbf{l}_1^N \in [d]^N, \forall N \in \mathbb{Z} \right\}$$

$$\mathbb{P}(\mathbf{y}_1^N = \mathbf{l}_1^N) = u' A^{(l_1)} A^{(l_1)} \dots A^{(l_N)} v, \quad \forall \mathbf{l}_1^N \in [d]^N$$

Realization problems

- ♦ **Informational complexity**
 - ✓ When is the minimal model identifiable from $\mathcal{S}^{(N)}$? window size N =?
- ♦ **Computational complexity**
 - ✓ Can we compute the minimal realization with efficient algorithms?
- ♦ **Statistical complexity**
 - ✓ Are the algorithms robust to estimation noise?

In general, learning HMMs is hard

Realization problems

- ♦ **Informational complexity**

- ✓ When is the minimal model identifiable from $\mathcal{S}^{(N)}$? window size N =?

Main contribution 1:

- ♦ **Computational complexity**

- ✓ Can we compute the minimal realization with efficient algorithms?

Main contribution 2

- ♦ **Statistical complexity**

- ✓ Are the algorithms robust to estimation noise?

In general, learning HMMs is hard



Quasi-HMM realization

Quasi-HMM realization

- ♦ **Input:** probabilities of length $N = 2n + 1$ strings

$$y_{-n}, \dots, y_{-1} \ y_0 \ y_1, \dots, y_n$$

$$\mathcal{S}^{(N)} = \{\mathbb{P}_{y_{-n}, \dots, y_n}\}$$

- ♦ **Output:** a minimal quasi-HMM realization

$$\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$$

- ♦ **Key:**

Arrange $\mathcal{S}^{(N)} = \{\mathbb{P}_{y_{-n}, \dots, y_n}\}$ into matrices:

$$H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$$

Matrix rank decomposition of $H^{(0)}$
equivalent minimal quasi realization up to linear transformation

Complexity?

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$
- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

Enumerate all length n strings: $L(\mathbf{l}_1^n) = (l_1 - 1)d^{n-1} + (l_2 - 1)d^{n-2} + \cdots + l_n, \quad \forall \mathbf{l}_1^n = [d]^n$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$\begin{aligned}[H^{(0)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} &= \mathbb{P}\left(\mathbf{y}_0^{n-1} = \mathbf{l}_1^n, \mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}\right) \\ &= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v\end{aligned}$$

n letters now and future

n letters in the past

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

Enumerate all length n strings: $L(\mathbf{l}_1^n) = (l_1 - 1)d^{n-1} + (l_2 - 1)d^{n-2} + \dots + l_n, \quad \forall \mathbf{l}_1^n = [d]^n$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$\begin{aligned}[H^{(0)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} &= \mathbb{P}\left(\mathbf{y}_0^{n-1} = \mathbf{l}_1^n, \mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}\right) \\ &= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v\end{aligned}$$

n letters now and future

n letters in the past

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

$$H^{(0)} = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} \begin{bmatrix} A^{(1)} \dots A^{(1)} v, \dots, A^{(d)} \dots A^{(d)} v \end{bmatrix}$$

Enumerate all length n strings: $L(\mathbf{l}_1^n) = (l_1 - 1)d^{n-1} + (l_2 - 1)d^{n-2} + \dots + l_n, \quad \forall \mathbf{l}_1^n = [d]^n$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$\begin{aligned}[H^{(0)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} &= \mathbb{P}\left(\mathbf{y}_0^{n-1} = \mathbf{l}_1^n, \mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}\right) \\ &= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v\end{aligned}$$

n letters now and future

n letters in the past

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

$$H^{(0)} = \begin{bmatrix} u' & & & (1) \\ u' & U & & (2) \\ u' & & & (d) \end{bmatrix} \begin{bmatrix} A^{(1)} \dots A & V' & A^{(d)} v \end{bmatrix}$$

Enumerate all length n strings: $L(\mathbf{l}_1^n) = (l_1 - 1)d^{n-1} + (l_2 - 1)d^{n-2} + \dots + l_n, \quad \forall \mathbf{l}_1^n = [d]^n$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$[H^{(j)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} = \mathbb{P}\left(\mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}, y_0 = j, \mathbf{y}_1^n = \mathbf{l}_1^n\right),$$

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$\begin{aligned}[H^{(j)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} &= \mathbb{P}\left(\mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}, y_0 = j, \mathbf{y}_1^n = \mathbf{l}_1^n\right), \\ &= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_0)}} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v\end{aligned}$$

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$[H^{(j)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} = \mathbb{P}\left(\mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}, y_0 = j, \mathbf{y}_1^n = \mathbf{l}_1^n\right),$$

$$= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_0)}} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v$$

n letters in the future
Conditional on current letter
n letters in the past

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

$$H^{(j)} = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} A^{(j)} \begin{bmatrix} A^{(1)} \dots A^{(1)} v, \dots, A^{(d)} \dots A^{(d)} v \end{bmatrix}$$

Quasi-HMM realization (algorithm)

- ♦ **Input:** $\mathcal{S}^{(N)}$ arranged into $H^{(0)}, \{H^{(j)} : j \in [d]\} \in \mathbb{R}^{d^n \times d^n}$

$$[H^{(j)}]_{L(\mathbf{l}_1^n), L(\mathbf{l}_{-1}^{-n})} = \mathbb{P}\left(\mathbf{y}_{-1}^{-n} = \mathbf{l}_{-1}^{-n}, y_0 = j, \mathbf{y}_1^n = \mathbf{l}_1^n\right),$$

$$= u' A^{(l_n)} \dots A^{(l_1)} \boxed{A^{(l_0)}} \boxed{A^{(l_{-1})} \dots A^{(l_{-n})}} v$$

n letters in the future
Conditional on current letter
n letters in the past

- ♦ assume minimal realization $\theta^o = (d, k, u, v, A^{(l)} : l \in [d])$

$$H^{(j)} = \begin{bmatrix} u' & & & & \\ u' & U & & & \\ u' & & & A^{(j)} & \\ & & & A^{(1)} & \cdot \\ & & & & V' \\ u' & & & & A^{(d)} \dots A^{(d)} v \end{bmatrix}$$

Quasi-HMM realization (algorithm)

$$d^n H^{(0)} =_{d^n} U \times k V'$$

matrix rank = k
If U, V are of full column rank

Algorithm 1 (Rank decomposition)

$$H^{(0)} = UV'$$

Quasi-HMM realization (algorithm)

$$d^n \quad \quad \quad k \quad \quad \quad d^n$$
$$d^n \quad H^{(0)} = d^n \quad U \quad \times \quad k \quad V'$$
$$H^{(j)} = \quad \quad \quad U \quad \times A^{(j)} \times \quad \quad \quad V'$$

Algorithm 1 (Rank decomposition)

$$H^{(0)} = UV'$$

$$\tilde{A}^{(j)} = U^\dagger H^{(j)} (V^\dagger)', \quad \forall j \in [d]$$

$$\tilde{u} = U' \mathbf{e}$$

$$\tilde{v} = V' \mathbf{e}$$

Unique up to linear transformation
Require U V have column rank = k

Quasi-HMM realization (informational complexity)

How large N needs to be so that U, V have full column rank = k?

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} T \in \mathbb{R}^{d^n \times k}$$

In practice, used as a heuristic algorithm

$$V = \begin{bmatrix} v' (A^{(1)} \dots A^{(1)})' \\ v' (A^{(1)} \dots A^{(2)})' \\ \vdots \\ v' (A^{(d)} \dots A^{(d)})' \end{bmatrix} (T^{-1})' \in \mathbb{R}^{d^n \times k}$$

Quasi-HMM realization (informational complexity)

How large N needs to be so that U, V have full column rank = k?

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} T \in \mathbb{R}^{d^n \times k}$$

$$V = \begin{bmatrix} v' (A^{(1)} \dots A^{(1)})' \\ v' (A^{(1)} \dots A^{(2)})' \\ \vdots \\ v' (A^{(d)} \dots A^{(d)})' \end{bmatrix} (T^{-1})' \in \mathbb{R}^{d^n \times k}$$

Best hope:

$$n \sim \mathcal{O}(\log_d k)$$

Quasi-HMM realization (informational complexity)

How large N needs to be so that U, V have full column rank = k?

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} T \in \mathbb{R}^{d^n \times k}$$

Theorem 1:
For **almost all** HMMs

$$V = \begin{bmatrix} v' (A^{(1)} \dots A^{(1)})' \\ v' (A^{(1)} \dots A^{(2)})' \\ \vdots \\ v' (A^{(d)} \dots A^{(d)})' \end{bmatrix} ($$

n > 2 \log_d(k)
guarantees U V full rank

Quasi-HMM realization (informational complexity)

$$\Theta_{(d,k)}^o = \{\theta^o : \text{order } k, \text{ alphabet size } d\} \quad N = 2n + 1$$

Theorem 1 (Information complexity for quasi-HMM realization).

- (1) Consider $\Theta_{(d,k)}^h$, the class of all HMMs with output alphabet size d and order k . There exists a measure zero set $\mathcal{E} \in \Theta_{(d,k)}^h$, such that for all the output process generated by HMMs in the set $\Theta_{(d,k)}^h \setminus \mathcal{E}$, the information in $\mathcal{S}^{(N)}$ is sufficient for computing the minimal quasi-HMM realization, for $N = 2n + 1$ with

$$n > 2 \log_d(k). \quad (16)$$

- (2) For any (d, k) pair, randomly generated a HMM in $\Theta_{(d,k)}^h$. If for a given window size $N = 2n + 1$, the matrix $H^{(0)} \in \mathbb{R}^{d^n \times d^n}$ constructed with $\mathcal{S}^{(N)}$ is of the maximal rank k , then for all output processes generated by $\Theta_{(d,k)}^o$, excluding a measure zero set, $\mathcal{S}^{(N)}$ is sufficient for computing the minimal quasi-HMM realization.

Quasi-HMM realization (informational complexity)

$$\Theta_{(d,k)}^o = \{\theta^o : \text{order } k, \text{ alphabet size } d\} \quad N = 2n + 1$$

Theorem 1 (Information complexity for quasi-HMM realization).

- (1) Consider $\Theta_{(d,k)}^h$, the class of all HMMs with output alphabet size d and order k . There exists a measure zero set $\mathcal{E} \in \Theta_{(d,k)}^h$, such that for all the output process generated by HMMs in the set $\Theta_{(d,k)}^h \setminus \mathcal{E}$, the information in $\mathcal{S}^{(N)}$ is sufficient for computing the minimal quasi-HMM realization, for $N = 2n + 1$ with

$$n > 2 \log_d(k). \quad (16)$$

- (2) For any (d, k) pair, randomly generated a HMM in $\Theta_{(d,k)}^h$. If for a given window size $N = 2n + 1$, the matrix $H^{(0)} \in \mathbb{R}^{d^n \times d^n}$ constructed with $\mathcal{S}^{(N)}$ is of the maximal rank k , then for all output processes generated by $\Theta_{(d,k)}^o$, excluding a measure zero set, is sufficient for computing the minimal quasi-HMM realization.

Polynomial computational complexity

$$\mathcal{O}(d(d^n)^3) \sim \mathcal{O}(dk^6)$$

Quasi-HMM realization (informational complexity)

$$\Theta_{(d,k)}^o = \{\theta^o : \text{order } k, \text{ alphabet size } d\} \quad N = 2n + 1$$

Theorem 1 (Information complexity for quasi-HMM realization).

- (1) Consider $\Theta_{(d,k)}^h$, the class of all HMMs with output alphabet size d and order k . There exists a measure zero set $\mathcal{E} \in \Theta_{(d,k)}^h$, such that for all the output process generated by HMMs in the set $\Theta_{(d,k)}^h \setminus \mathcal{E}$, the information in $\mathcal{S}^{(N)}$ is sufficient for computing the minimal quasi-HMM realization, for $N = 2n + 1$ with

$$n > 2 \log_d(k). \quad (16)$$

- (2) For any (d, k) pair, randomly generated a HMM in $\Theta_{(d,k)}^h$. If for a given window size $N = 2n + 1$, the matrix $H^{(0)} \in \mathbb{R}^{d^n \times d^n}$ constructed with $\mathcal{S}^{(N)}$ is of the maximal rank k , then for all output processes generated by $\Theta_{(d,k)}^o$, excluding a measure zero set, is sufficient for computing the minimal quasi-HMM realization.

Polynomial computational complexity

$$\mathcal{O}(d(d^n)^3) \sim \mathcal{O}(dk^6)$$

- ♦ For almost all HMM, min order quasi realization is easy
- ♦ However, there exist information theoretic hard cases

Sketch of the proof

- We want to show for generic choices of Q and O , if $N = 2n + 1$ satisfies $n > 2 \log_d(k)$, the matrices U, V have full column rank k .
- Since the minors of U and V are nonzero polynomials in the elements of Q and O , it is enough to show for some specific choice of Q and O .
- $Q \in \mathbb{R}^{k \times k}$: the state shifting matrix

$$Q_{i-1,i} = 1, \text{ for } 2 \leq i \leq k, \text{ and } Q_{k,1} = 1,$$

$O \in \mathbb{R}^{d \times k}$: columns are independent unit-norm d -dimensional random vectors.

The i -th column of U is given by

$$U_{[:,i]} = O_{[:,i]} \odot \dots \odot O_{[:,i+n-1]}.$$

- Show that with probability greater than 0 (random matrix, concentration ineq)

$$\sigma_{\min}(U'U) > 0.$$

Minimal quasi-HMM realization (summary)

- ♦ Key to quasi-HMM realization:
 - ✓ conditional independence,
 - ✓ matrix decomposition, unique up to linear transformation
 - ✓ Minimal order is determined by the rank of $H^{(0)}$
- ♦ For **almost all** HMM, minimal quasi-HMM realization is easy
 - ✓ Informational : $N \sim \log_d(k)$
 - ✓ Computational : matrix factorization $\mathcal{O}(dk^6)$

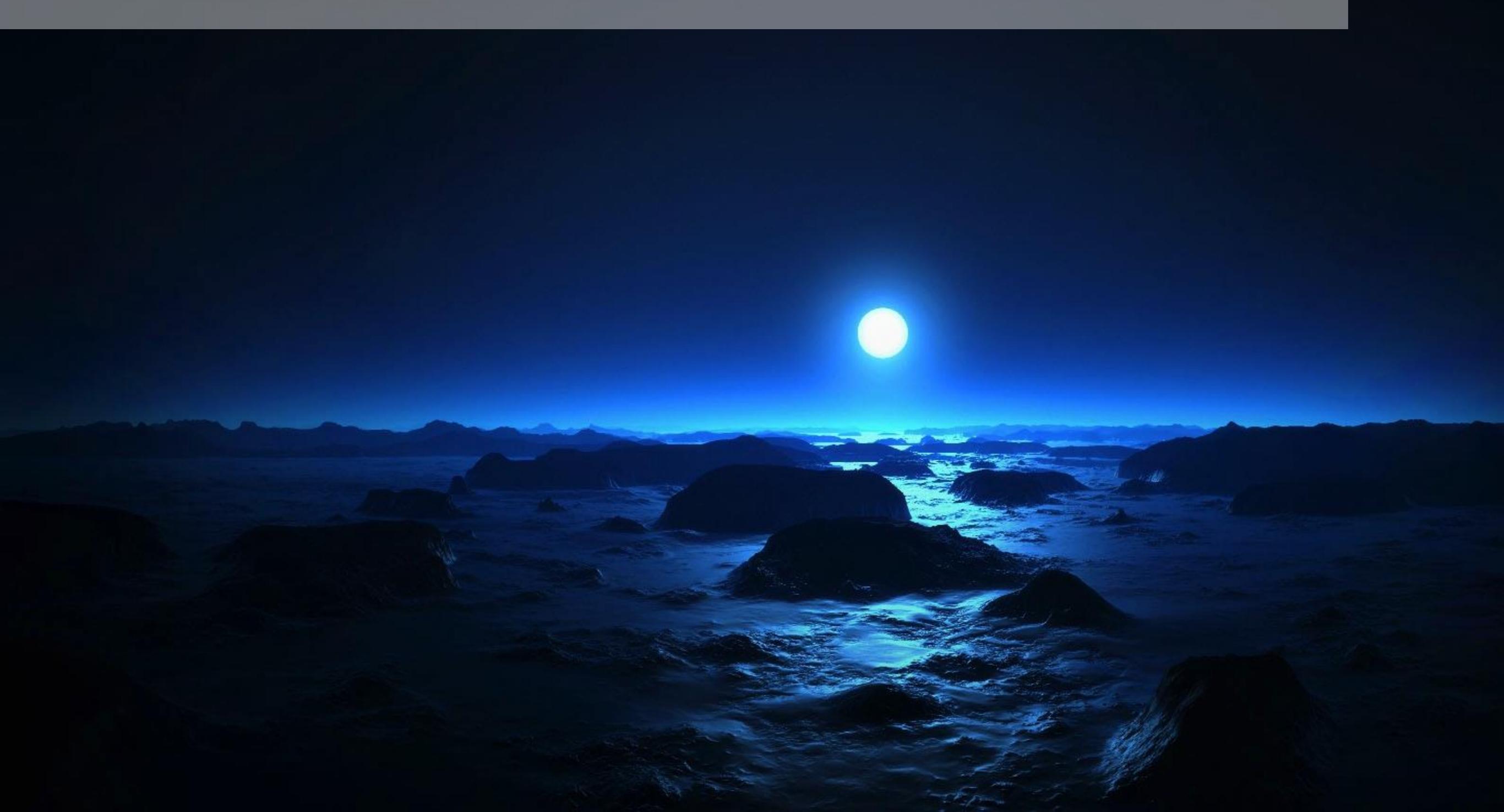
Minimal quasi-HMM realization (summary)

- ♦ Key to quasi-HMM realization:
 - ✓ conditional independence,
 - ✓ matrix decomposition, unique up to linear transformation
 - ✓ Minimal order is determined by the rank of $H^{(0)}$
- ♦ For **almost all** HMM, minimal quasi-HMM realization is easy
 - ✓ Informational : $N \sim \log_d(k)$
 - ✓ Computational : matrix factorization $\mathcal{O}(dk^6)$

Minimal quasi-HMM realization (summary)

- ♦ Key to quasi-HMM realization:
 - ✓ conditional independence,
 - ✓ matrix decomposition, unique up to linear transformation
 - ✓ Minimal order is determined by the rank of $H^{(0)}$
- ♦ For **almost all** HMM, minimal quasi-HMM realization is easy
 - ✓ Informational : $N \sim \log_d(k)$
 - ✓ Computational : matrix factorization $\mathcal{O}(dk^6)$
- ✓ Smoothed analysis VS generic analysis

Minimal HMM realization

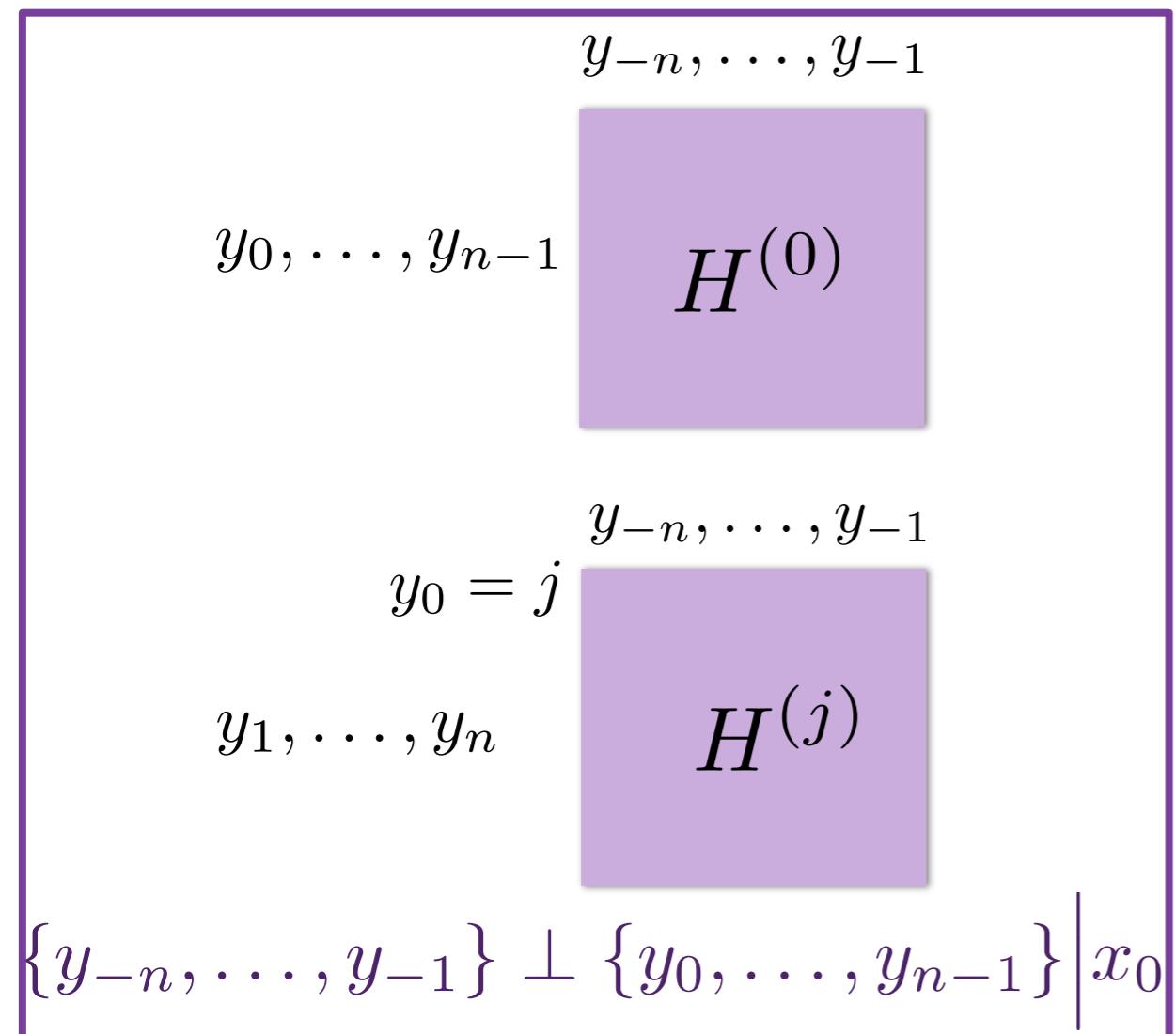


Minimal HMM realization

- ♦ HMM conditional independence : probabilities in product form

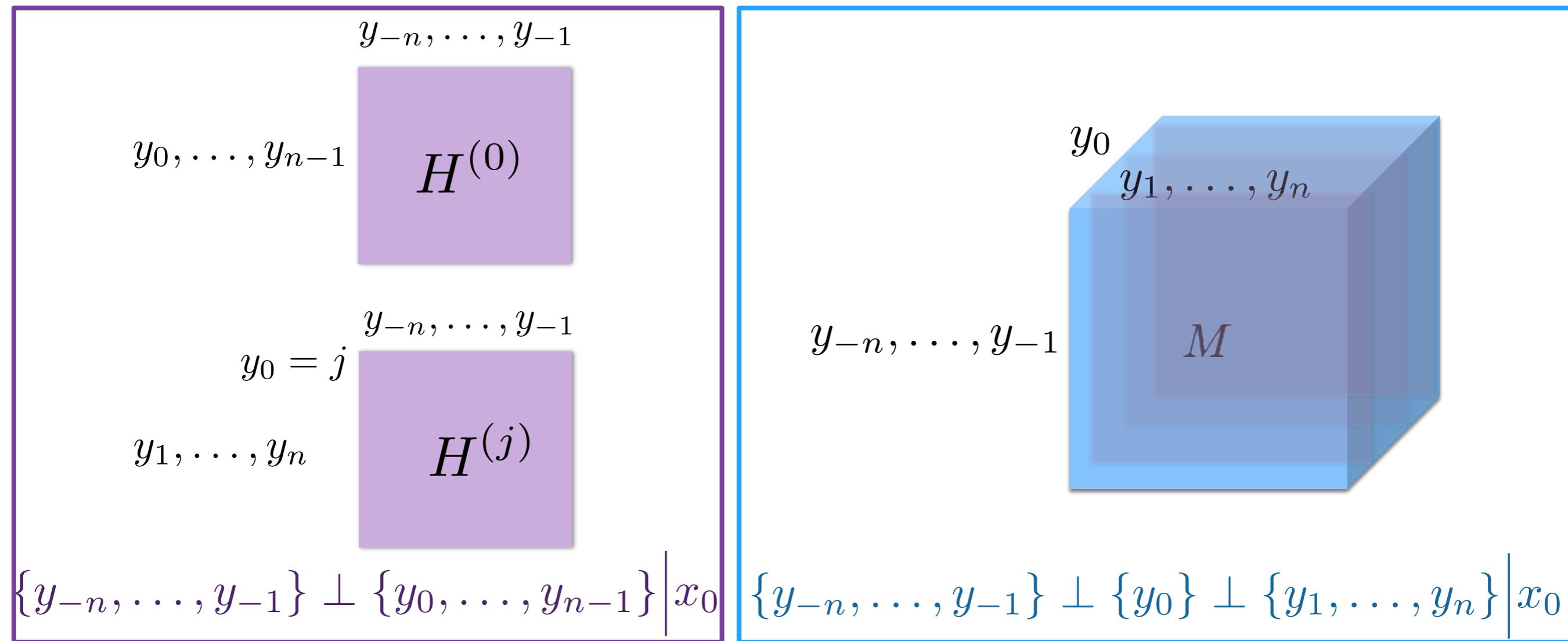
Minimal HMM realization

- ♦ HMM conditional independence : probabilities in product form



Minimal HMM realization

- ♦ HMM conditional independence : probabilities in product form
- ♦ Exploit the uniqueness of tensor decomposition



Minimal HMM realization

- ◆ **Input:** $\mathcal{S}^{(N)} = \{\mathbb{P}_{y_{-n}, \dots, y_n}\} \quad N = 2n + 1$
- ◆ **Output:** a minimal HMM model $\theta^h = (d, k, Q, O)$
- ◆ **Key:**

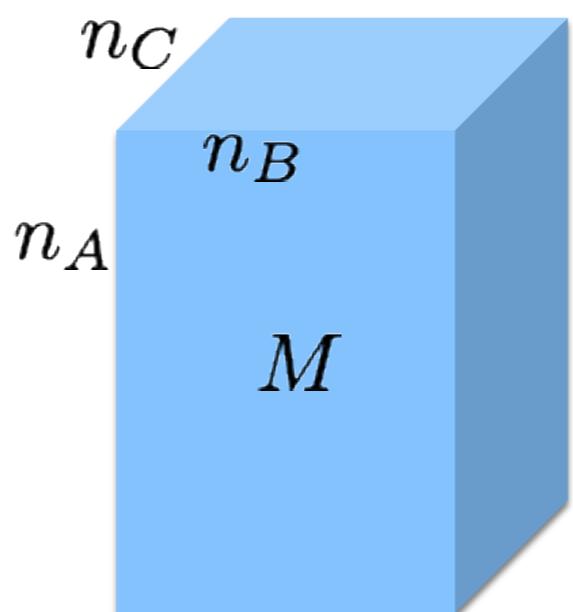
Arrange $\mathcal{S}^{(N)}$ into a three way tensor

$$M \in \mathbb{R}^{d^n \times d^n \times d}$$

Tensor rank decomposition
recover minimal HMM realization up to state relabeling

Preliminaries on Tensor

- ♦ Tensor definition
 - ✓ Multi-way array
 - ✓ Multi-linear mapping



$$M_{j_1, j_2, j_3}, \quad j_1 \in [n_A], j_2 \in [n_B], j_3 \in [n_C]$$

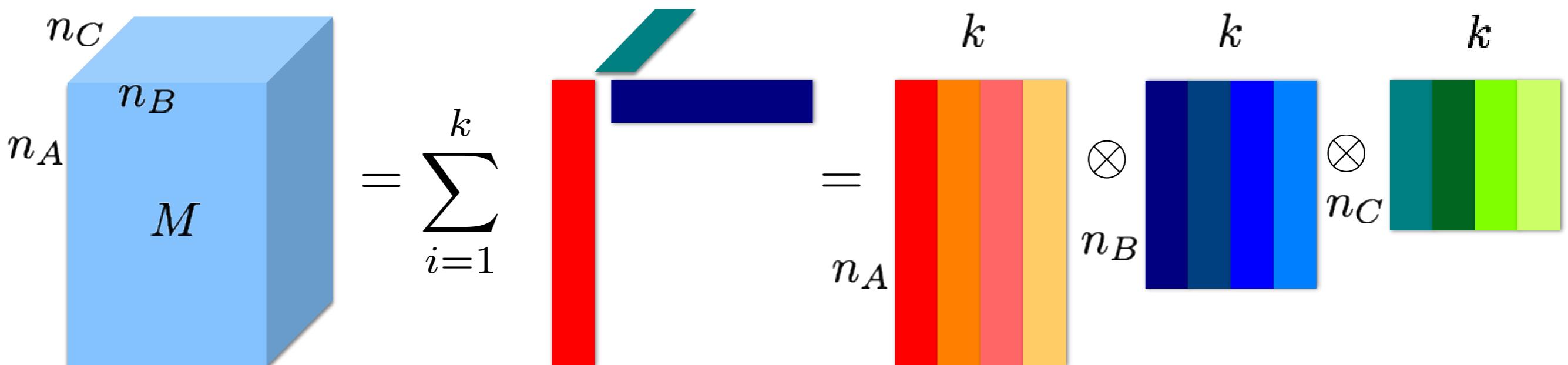
Compare to matrix H_{j_1, j_2}

Preliminaries on Tensor

- ♦ Tensor rank decomposition (CP/PARAFAC)

- ✓ Sum of rank one tensors $[\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}]_{j_1,j_2,j_3} = a_{j_1} b_{j_2} c_{j_3}$

$$M = \sum_{i=1}^k A_{[:,i]} \otimes B_{[:,i]} \otimes C_{[:,i]} = A \quad \otimes \quad B \quad \otimes \quad C$$



Compare to matrix rank decomposition $H = U \otimes V = \sum_i^k U_{[:,i]} V_{[:,i]}^\top$

Minimal HMM realization (intuition)

- ♦ **Key:** $S^{(N)}$ arranged into a three way tensor $M \in \mathbb{R}^{d^n \times d^n \times d}$

$$M_{L(\mathbf{l}_{-1}^{-n}), L(\mathbf{l}_1^n), l_0} = \mathbb{P}(\mathbf{y}_{-n}^n = \mathbf{l}_{-n}^n)$$

M relates to the model parameters via unique tensor factorization

$$M = A \otimes B \otimes C$$

$$A, B \in \mathbb{R}^{d^n \times k} \text{ and } C \in \mathbb{R}^{d \times k}$$

$$A = \mathbb{P}(y_1, y_2, \dots, y_n \mid x_0)$$

$$B = \mathbb{P}(y_{-1}, y_{-2}, \dots, y_{-n} \mid x_0)$$

$$C = \mathbb{P}(y_0, x_0)$$

Minimal HMM realization (intuition)

- ♦ **Key:** $S^{(N)}$ arranged into a three way tensor $M \in \mathbb{R}^{d^n \times d^n \times d}$

$$M_{L(\mathbf{l}_{-1}^{-n}), L(\mathbf{l}_1^n), l_0} = \mathbb{P}(\mathbf{y}_{-n}^n = \mathbf{l}_{-n}^n)$$

M relates to the model parameters via unique tensor factorization

$$M = A \otimes B \otimes C$$

$$A, B \in \mathbb{R}^{d^n \times k} \text{ and } C \in \mathbb{R}^{d \times k}$$

$$A = \mathbb{P}(y_1, y_2, \dots, y_n \mid x_0)$$

$$A = \underbrace{(O \odot (O \odot (\underbrace{O \odot \dots (O \odot O}_n Q) \dots)_Q)_Q)}_n Q,$$

$$B = \mathbb{P}(y_{-1}, y_{-2}, \dots, y_{-n} \mid x_0)$$

$$B = \underbrace{(O \odot (O \odot (\underbrace{O \odot \dots (O \odot O \tilde{Q}) \dots}_n \tilde{Q}) \dots)_\tilde{Q}) \tilde{Q}}_n \tilde{Q},$$

$$C = \mathbb{P}(y_0, x_0)$$

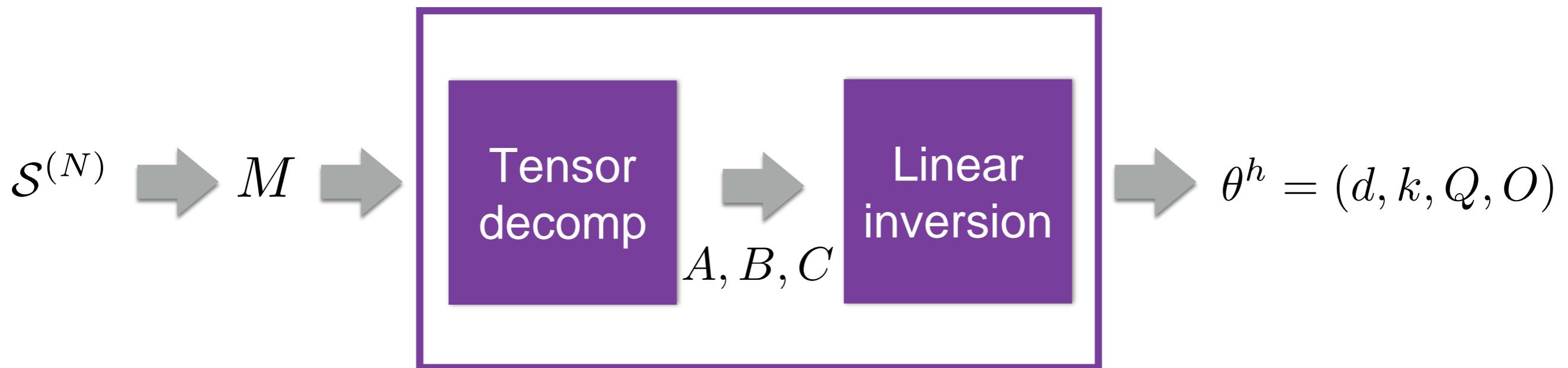
$$C = Odiag(\pi)$$

Minimal HMM realization (two-step approach)

$$\begin{matrix} d \\ d^n \\ M \end{matrix} = \begin{matrix} k \\ d^n \\ A \end{matrix} \otimes \begin{matrix} k \\ d^n \\ B \end{matrix} \otimes \begin{matrix} k \\ d \\ C \end{matrix} \Leftrightarrow Q, O \end{matrix}$$

A diagram illustrating the minimal HMM realization. On the left, a blue cube labeled M has dimensions d (depth), d^n (height), and d^n (width). An equals sign follows, followed by three tensors. The first tensor has dimensions k (depth), d^n (height), and d^n (width), with the label A in its middle section. The second tensor has dimensions k (depth), d^n (height), and d^n (width), with the label B in its middle section. The third tensor has dimensions k (depth), d (height), and k (width), with the label C in its middle section. Between the first and second tensors is a circle with a cross symbol (\otimes). Between the second and third tensors is another circle with a cross symbol (\otimes). To the right of the third tensor is a double-headed arrow (\Leftrightarrow) followed by Q, O .

Two-step realization algorithm



Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) M \Leftrightarrow A, B, C \rightarrow (2) A, B, C \Leftrightarrow (Q, O)$$

Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) \quad M \Leftrightarrow A, B, C \rightarrow (2) \quad A, B, C \Leftrightarrow (Q, O)$$

Proposition 2: sufficient conditions for invertibility of (1)

$$\text{krank}(A) + \text{krank}(B) + \text{krank}(C) \geq 2k + 2,$$

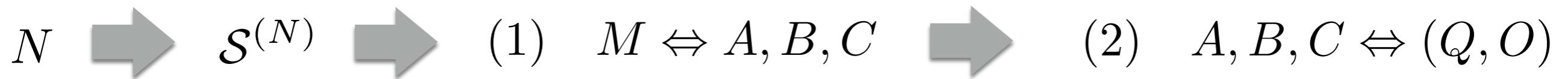
Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) \quad M \Leftrightarrow A, B, C \rightarrow (2) \quad A, B, C \Leftrightarrow (Q, O)$$

Theorem 3: sufficient conditions for invertibility of (1)

$$\text{rank}(A) = \text{rank}(B) = k$$

Minimal HMM realization (identifiability)



Theorem 3: sufficient conditions for invertibility of (1)

$$\text{rank}(A) = \text{rank}(B) = k$$

Theorem:

For almost all HMMs

$$n > 2 \log_d(k)$$

guarantees A,B full rank

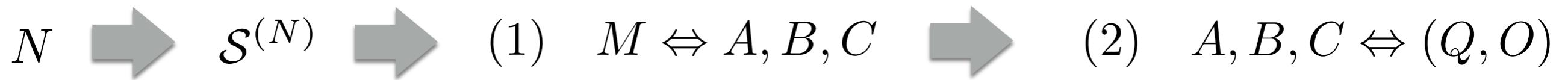
Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) \quad M \Leftrightarrow A, B, C \rightarrow (2) \quad A, B, C \Leftrightarrow (Q, O)$$

Theorem 2: sufficient conditions for invertability of (2)

1. If A has full column rank
2. Or if O has full column rank

Minimal HMM realization (algorithm)



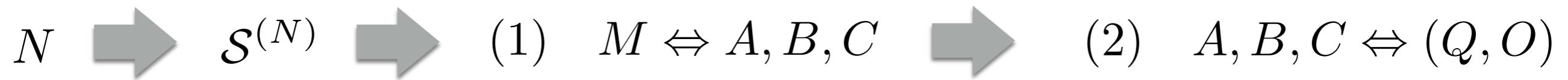
Theorem 2: algorithm for inverting (2)

Obtain the observation probabilities by: $O_{[:,i]} = C_{[:,i]} / (\mathbf{e}' C_{[:,i]}), \quad \forall i \in [k]$. Given the matrix A , i.e., $\mathbb{P}(\mathbf{y}_1^n | x_0)$, we can marginalize to obtain the matrices:

$$\tilde{A}^{(1)} = \mathbb{P}(y_1 | x_0) \quad \tilde{A}^{(n-1)} = \mathbb{P}(\mathbf{y}[1 : n-1] | x_0),$$

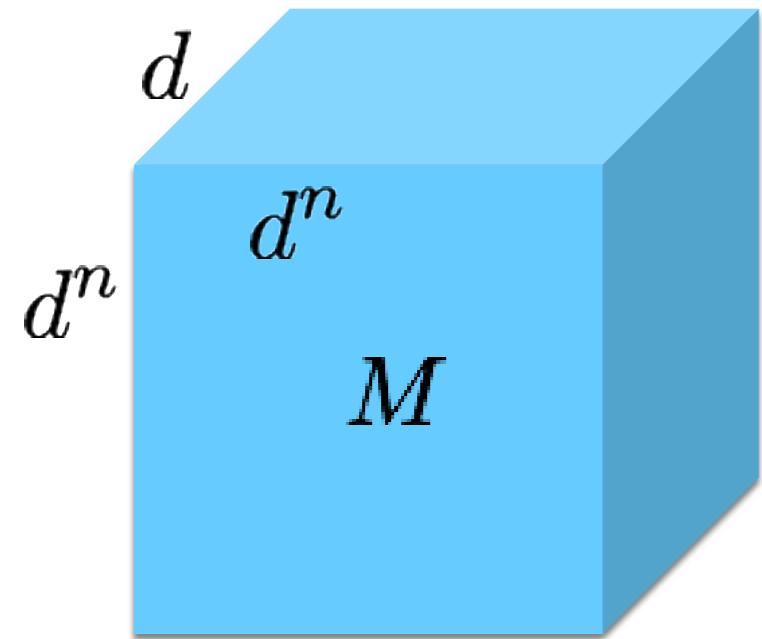
1. If A has full column rank: $Q = (O \odot \tilde{A}^{(n-1)})^{-1} A'$;
2. if O has full column rank: $Q = O^\dagger \tilde{A}^{(1)}$.

Minimal HMM realization (algorithm)



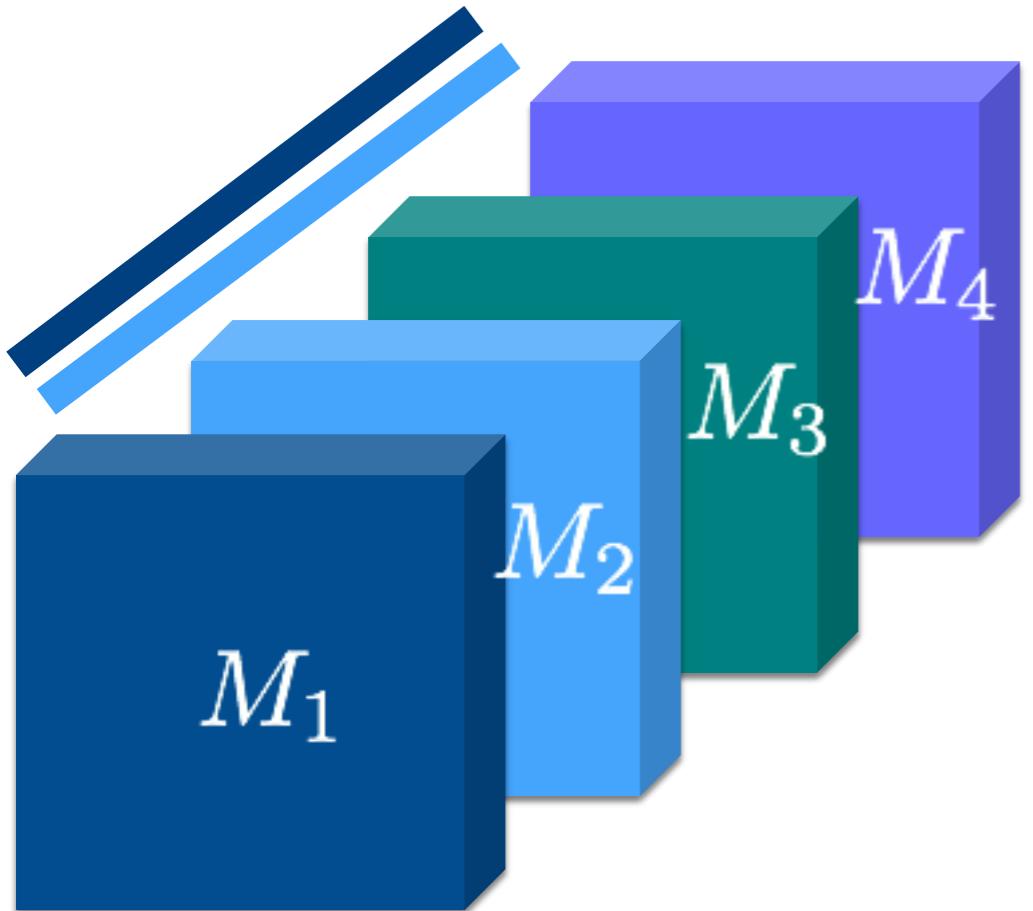
- ♦ Unique tensor decomposition to invert (1)
 - ✓ In general tensor decomposition is hard
(Alternating least square)
 - ♦ Have efficient algorithm if
 - ✓ A and B have full column rank ← Algorithm 2
 - ✓ C is full rank, $A \odot B$ is full rank and... ← Algorithm 3
 - ✓ The two cases apply to *almost all* HMMs for large enough N

Algorithm 2 (Simultaneous diagonalization)

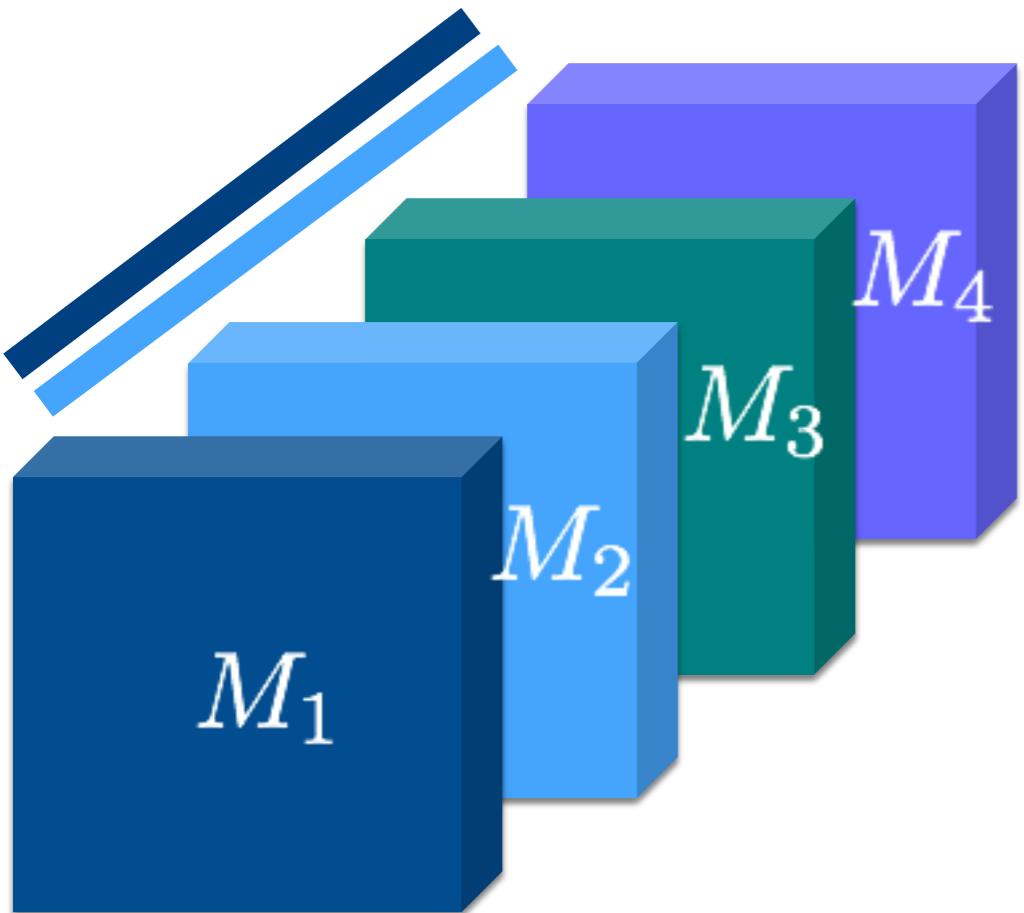


- ◆ Take 2 random projections along the 3rd dim

Algorithm 2 (Simultaneous diagonalization)



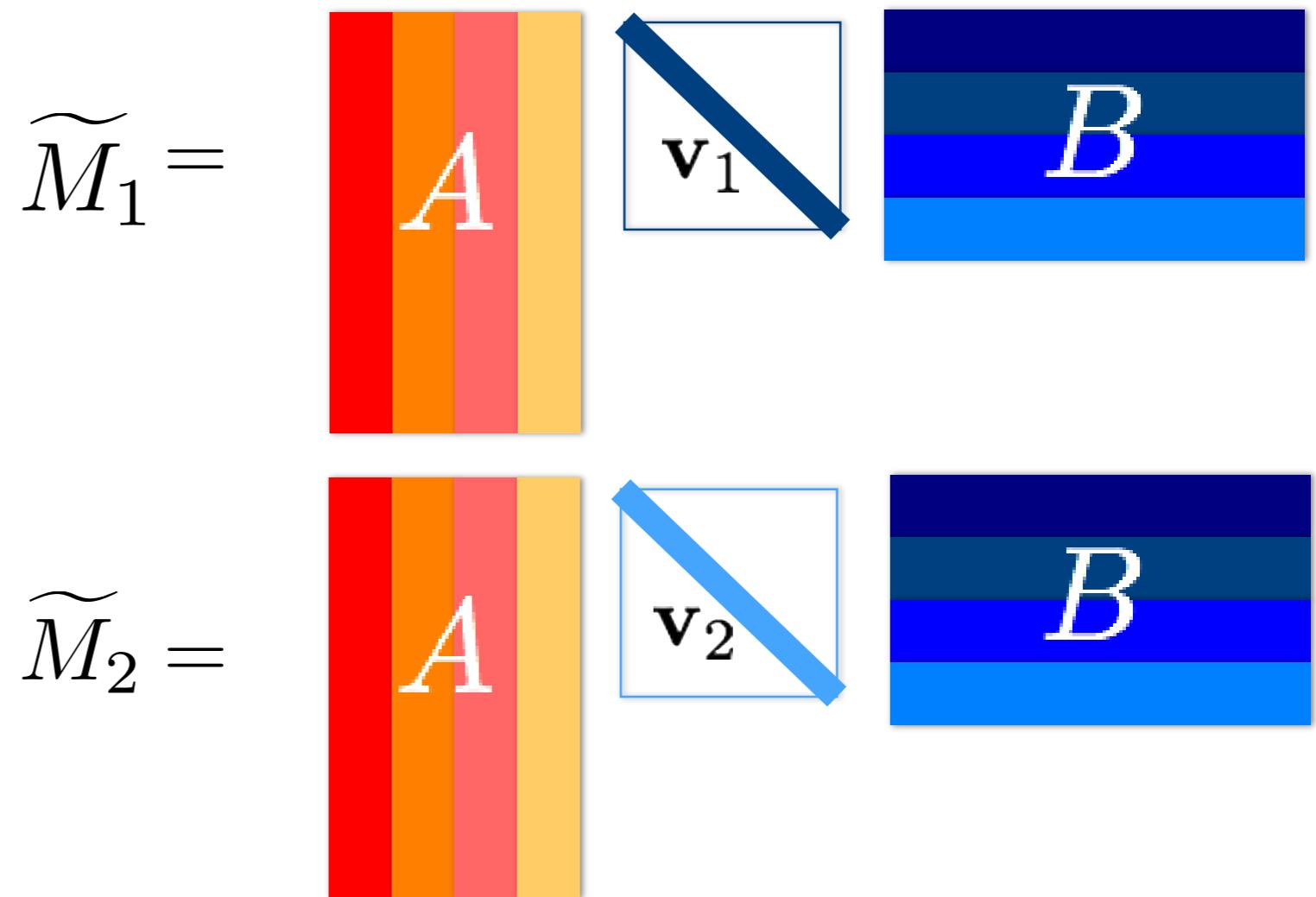
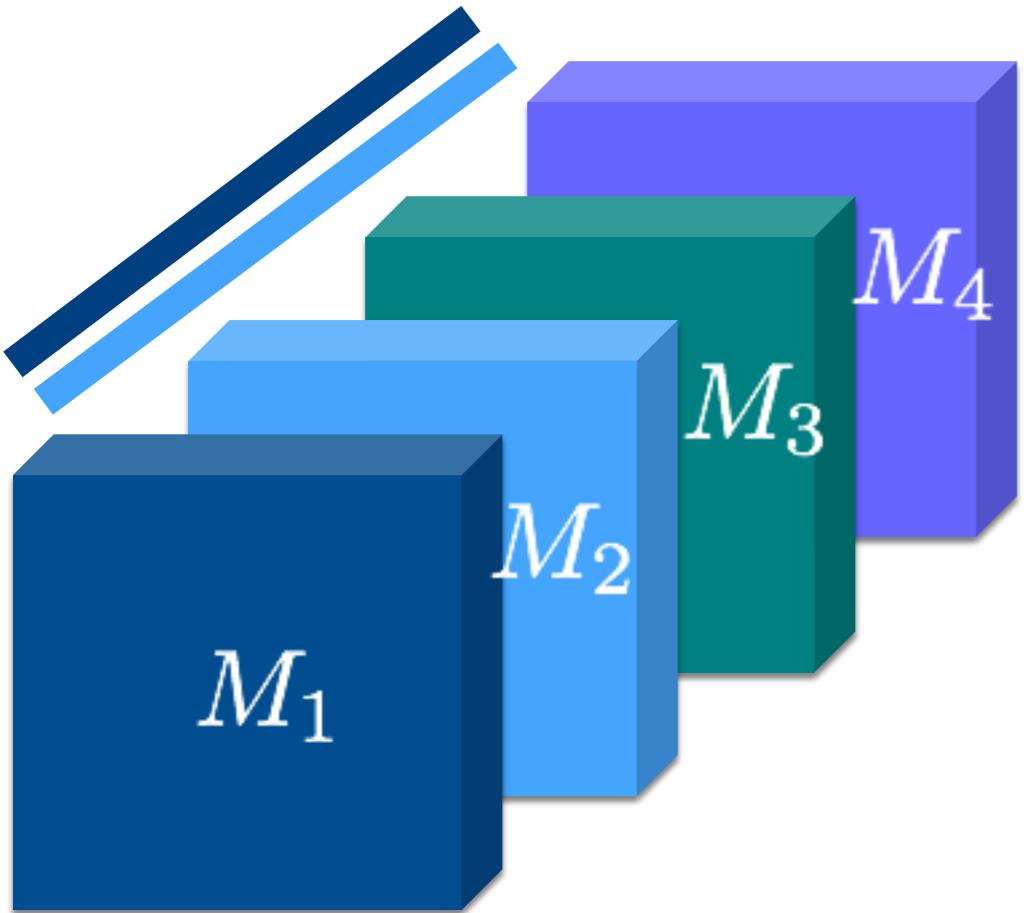
Algorithm 2 (Simultaneous diagonalization)



$$\widetilde{M}_1 = \begin{array}{c} \text{red slice} \\ \text{orange slice} \\ A \\ \text{pink slice} \\ \text{yellow slice} \end{array}$$
$$\mathbf{v}_1$$
$$B$$
$$\widetilde{M}_2 = \begin{array}{c} \text{red slice} \\ \text{orange slice} \\ A \\ \text{pink slice} \\ \text{yellow slice} \end{array}$$
$$\mathbf{v}_2$$
$$B$$

- ◆ Take 2 random projections along the 3rd dimension → 2 matrix slices
- ◆ Simultaneous diagonalization is unique almost surely

Algorithm 2 (Simultaneous diagonalization)



- ♦ Take 2 random projections along the 3rd dimension

$$\widetilde{M}_1 = M(I_{d^n \times d^n}, I_{d^n \times d^n}, \mathbf{v}_1), \quad \widetilde{M}_2 = M(I_{d^n \times d^n}, I_{d^n \times d^n}, \mathbf{v}_2).$$

$$\begin{aligned} \widetilde{M}_1 \widetilde{M}_2^{-1} &= A \Lambda A^{-1}, \\ \widetilde{M}_2 \widetilde{M}_1^{-1} &= B \Lambda^{-1} B^{-1}. \end{aligned} \Rightarrow \text{pair up eigenvalues to get } A, B$$

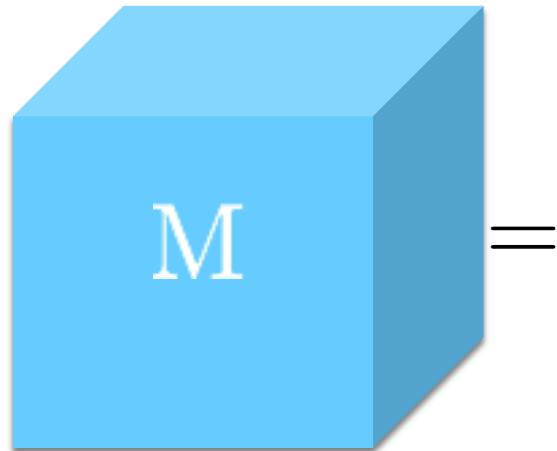
solve linear eqn to obtain C

Algorithm 3 (FOOBI)

$$(1) \quad M = A \otimes B \otimes C$$

A,B are rank degenerate
But if $A \odot B$ has full column rank

- ♦ Degenerate case: Q is not full rank, set N = 3

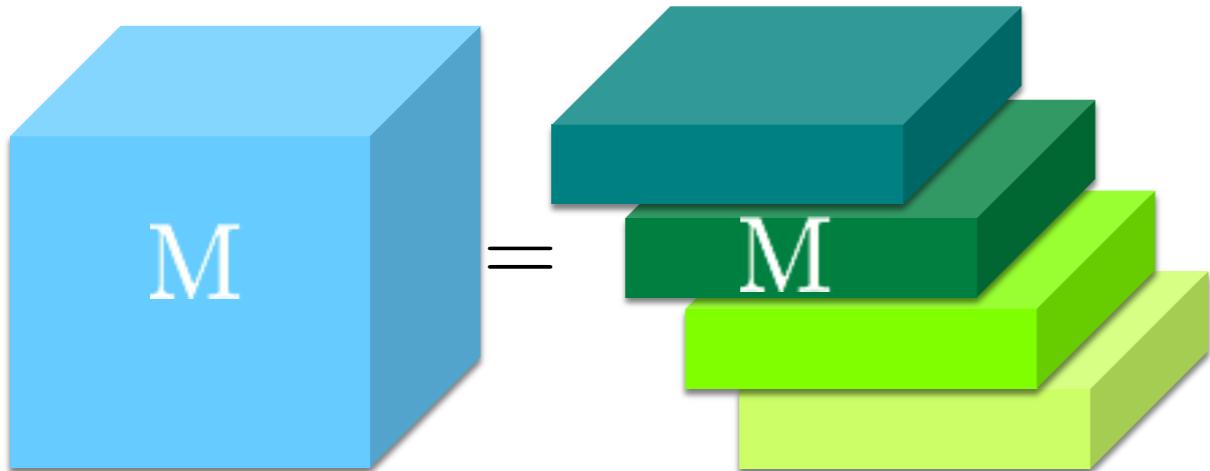


Algorithm 3 (FOOBI)

$$(1) \quad M = A \otimes B \otimes C$$

A,B are rank degenerate
But if $A \odot B$ has full column rank

- ♦ Degenerate case: Q is not full rank, set N = 3
- ♦ Matricization M along the 3rd dim: \overline{M}

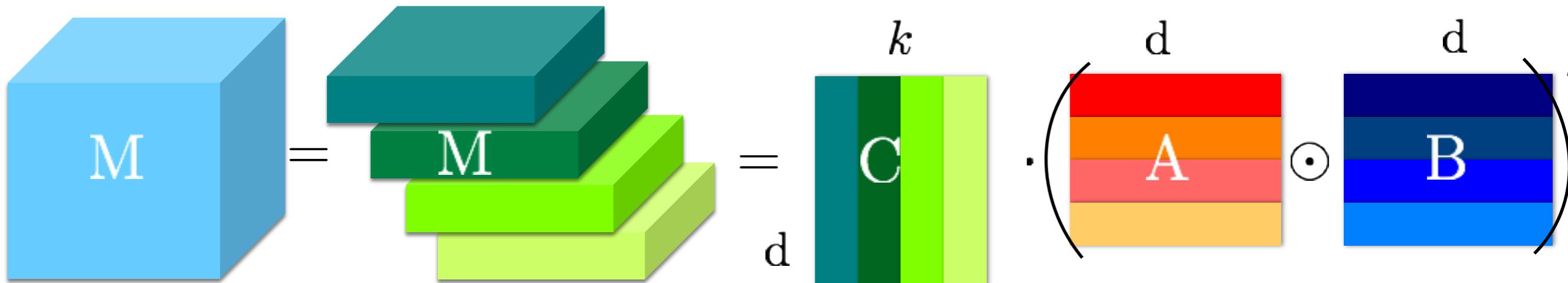


Algorithm 3 (FOOBI)

$$(1) \quad M = A \otimes B \otimes C$$

A,B are rank degenerate
But if $A \odot B$ has full column rank

- ♦ Degenerate case: Q is not full rank, set N = 3
- ♦ Matricization M along the 3rd dim: \overline{M}



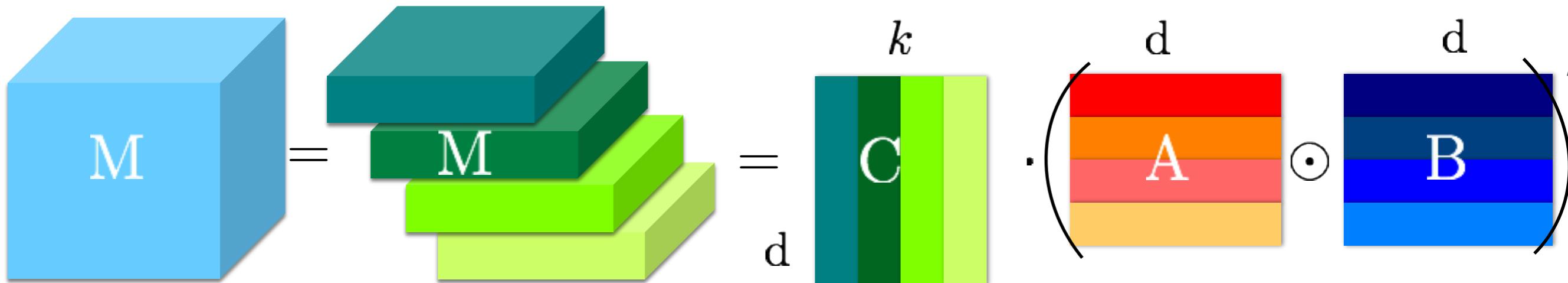
$$\overline{M} = C \otimes (A \odot B)'$$

Algorithm 3 (FOOBI)

$$(1) \quad M = A \otimes B \otimes C$$

A,B are rank degenerate
But if $A \odot B$ has full column rank

- ♦ Degenerate case: Q is not full rank, set N = 3
- ♦ Matricization M along the 3rd dim: \overline{M}



$$\overline{M} = C \otimes (A \odot B)'$$

Can be converted to
solving linear equations and
simultaneous diagonalization

rank decomposition $\overline{M} = UV'$

such that: each column of V is a rank 1 matrix

Algorithm 3 (correctness)

$$\Theta_{(d,k,r)}^h = \{\theta^h : \text{order } k, \text{ alphabet size } d, Q \text{ is of rank } r < k\}$$

Theorem 4: Correctness of Algorithm 3

Pick a random instance from the model class $\Theta_{(d,k,r)}^h$

If Algorithm 3 works correctly to recover the minimal HMM

Then it works for almost all instances in $\Theta_{(d,k,r)}^h$

Minimal HMM realization (summary)

- ♦ Key to HMM realization:
 - ✓ conditional independence: **current, past** and **future**
 - ✓ tensor decomposition, unique up to column permutation
- ♦ For almost all HMM, minimal HMM realization is easy
 - ✓ Informational : $N \sim \log_d(k)$
 - ✓ Computational : simultaneous diagonalization $\mathcal{O}(dk^6)$
- ♦ A class of degenerate cases can also be efficiently realized

Minimal HMM realization (summary)

- ♦ Key to HMM realization:
 - ✓ conditional independence: **current, past and future**
 - ✓ tensor decomposition, unique up to column permutation
- ♦ For almost all HMM, minimal HMM realization is easy
 - ✓ Informational : $N \sim \log_d(k)$
 - ✓ Computational : simultaneous diagonalization $\mathcal{O}(dk^6)$
- ♦ A class of degenerate cases can also be efficiently realized

Summary & Discussions



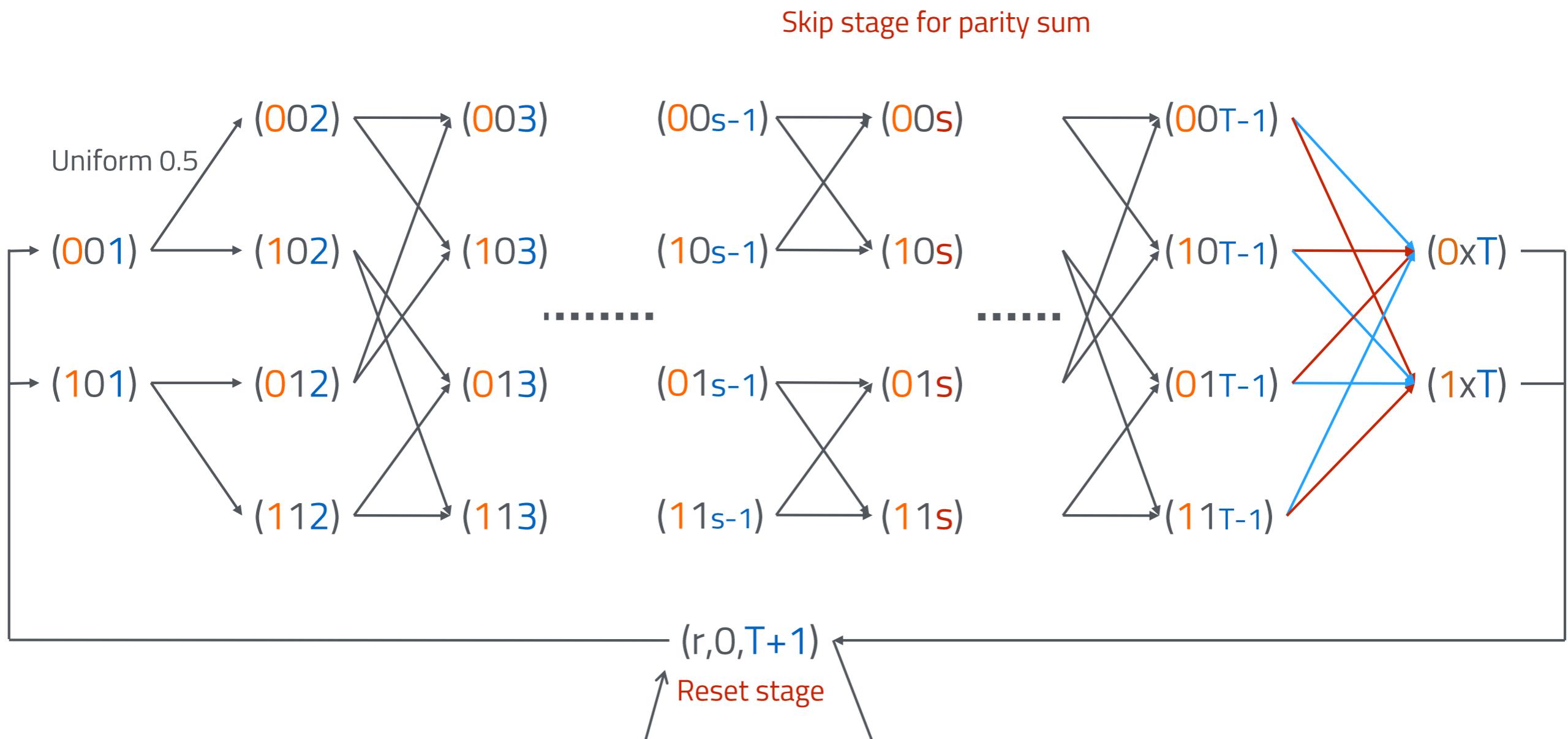
Discussions

- ♦ Hardness of worst cases
 - ✓ Reduction to parity of noise
- ♦ Learning via realization algorithms
 - ✓ Statistical complexity
- ♦ Ongoing works
 - ✓ Smoothed analysis
 - ✓ Agnostic learning
 - ✓ HMM model reduction

Hardness results

- ♦ We showed for almost all HMM, operator model and HMM model both are easy to learn (poly time algorithm, poly sample complexity)
- ♦ There exist information theoretic **hard** cases [Abe,Warmuth] [Kearns]
 - ✓ Given order k , cannot learn $\theta^h = (d, k, Q, O)$ in poly time, unless RP=NP
 - ✓ Unknown k , HMM is not efficiently PAC learnable, under noisy parity assumption
- ♦ Example: reduction to parity of noise
 - ✓ Number of states = $4T-1$, observation alphabet = 2
 - ✓ required window size = T
 - ✓ There exists an η such that there is no efficient algorithm for learning parity under uniform distribution in the PAC framework with classification noise η

State = (Emission, Parity sum, Stage)



HMM learning problems

- ♦ Observe: sample sequences of the output process of an HMM
- ♦ Goal: an operator / HMM model to fit the process

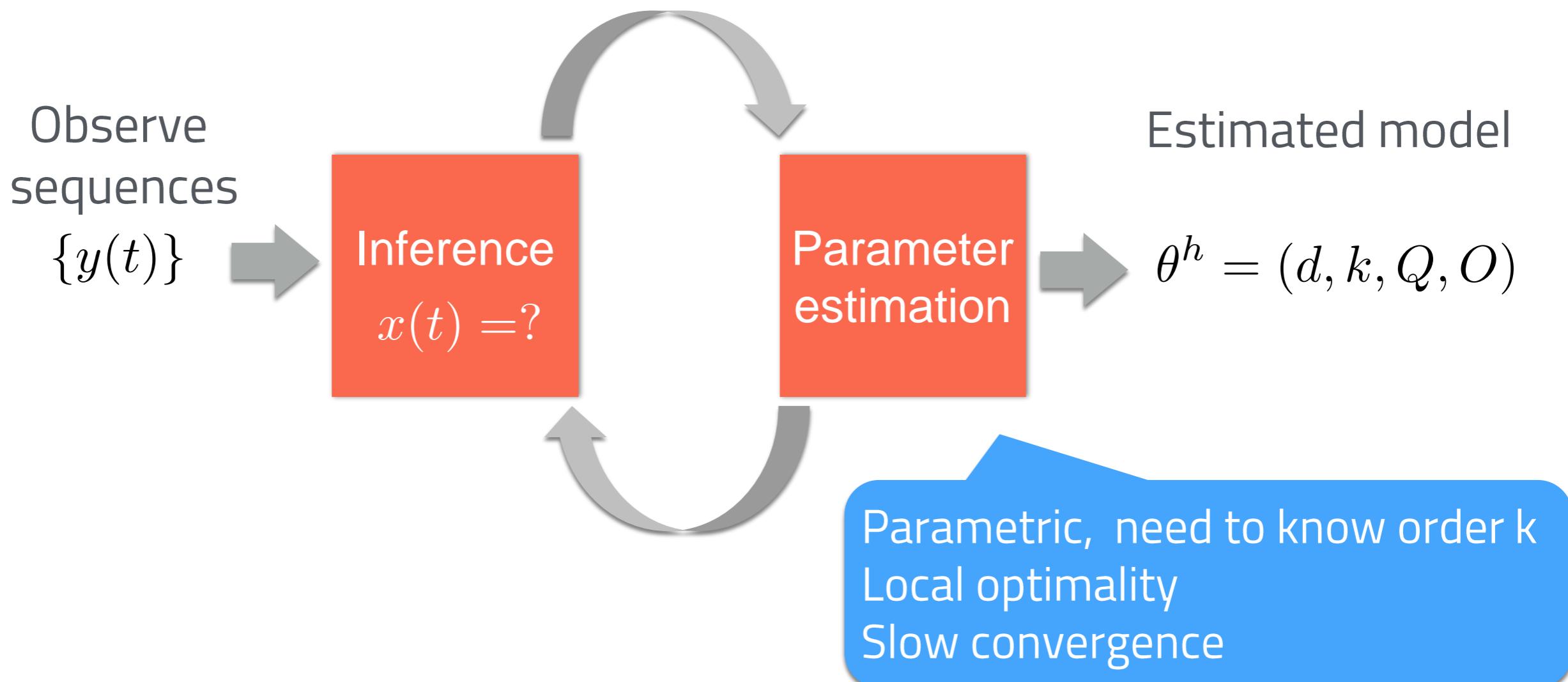
Maximal likelihood
estimator

- ♦ Old art: EM (Expectation-Maximization)
- ♦ New art: Spectral method (subspace id/ automata / tensor)

Statistics matching

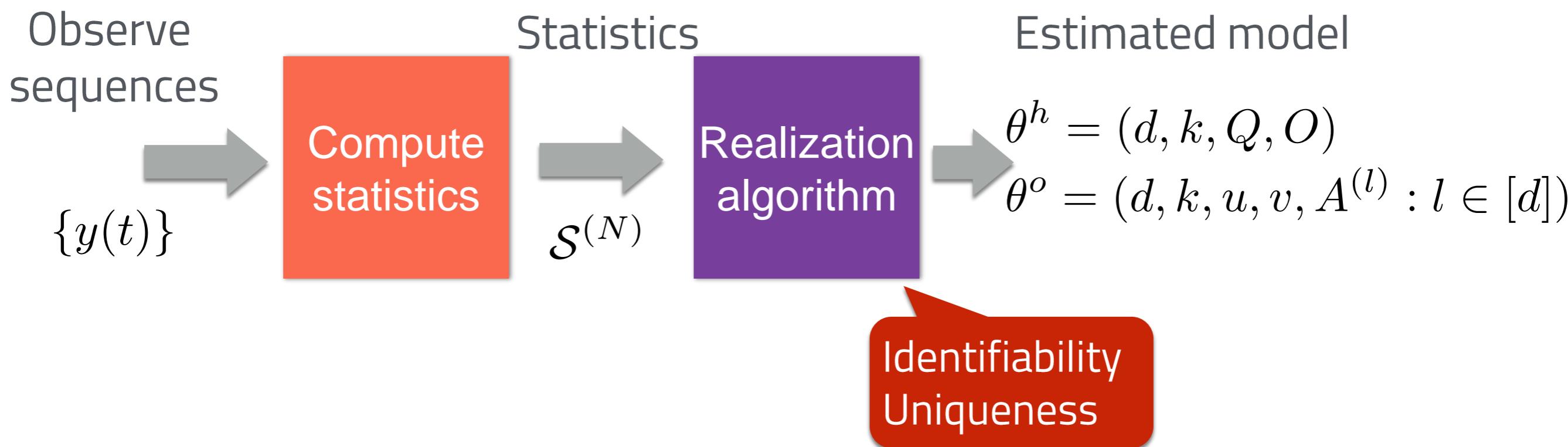
HMM learning problems

- ♦ Observe: sample sequences of the output process of an HMM
- ♦ Goal: an operator / HMM model to fit the process
- ♦ **Old art:** EM



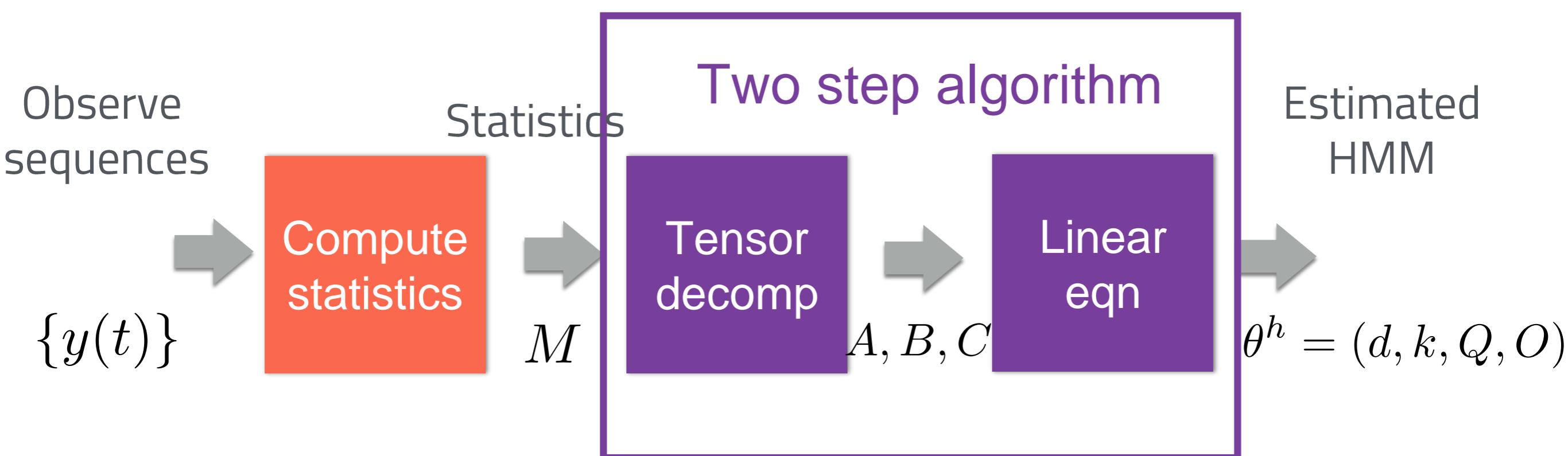
HMM learning problems

- ♦ Observe: sample sequences of the output process of an HMM
- ♦ Goal: an operator / HMM model to fit the process
- ♦ **New art:** spectral method



MinimalHMM learning

- ♦ Observe: sample sequences of the output process of an HMM
- ♦ Goal: an operator / HMM model to fit the process
- ♦ **New art:** spectral method



Sample complexity

Ongoing works

- ♦ Smoothed analysis – perturbation based
Understand the polynomial time algorithm in practice
- ♦ Agnostic HMM learning
If the underlying model is not a HMM...
which class of models have better fitting power?
- ♦ HMM model reduction

Appendix

Quasi-HMM realization (correctness of the algorithm)

Proposition 1. Let $\theta^o = (d, k, u, v, A^{(j)} : j \in [d])$ be a minimal order quasi realization for the output process. If the matrix $H^{(0)}$ has rank k , then the Algorithm returns an operator model $\tilde{\theta}^o$ that is equivalent to θ^o up to linear transformation.

Assume the output process can be realized by an HMM model of order k^h . The order k^h is always lower bounded by the order of the minimal order quasi realization k .

- ◆ Sketch of the proof

Quasi-HMM realization (correctness of the algorithm)

Proposition 1. Let $\theta^o = (d, k, u, v, A^{(j)} : j \in [d])$ be a minimal order quasi realization for the output process. If the matrix $H^{(0)}$ has rank k , then the Algorithm returns an operator model $\tilde{\theta}^o$ that is equivalent to θ^o up to linear transformation.

Assume the output process can be realized by an HMM model of order k^h . The order k^h is always lower bounded by the order of the minimal order quasi realization k .

- ♦ Sketch of the proof

$$H^{(0)} = UV', \quad U, V \text{ both of full rank } k$$

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} \quad T = \begin{bmatrix} \tilde{u}' \tilde{A}^{(1)} \dots \tilde{A}^{(1)} \\ \tilde{u}' \tilde{A}^{(1)} \dots \tilde{A}^{(2)} \\ \vdots \\ \tilde{u}' \tilde{A}^{(d)} \dots \tilde{A}^{(d)} \end{bmatrix}$$

$$V' = T^{-1} [A^{(1)} \dots A^{(1)} v, \dots, A^{(d)} \dots A^{(d)} v] = [\tilde{A}^{(1)} \dots \tilde{A}^{(1)} v, \dots, \tilde{A}^{(d)} \dots \tilde{A}^{(d)} \tilde{v}]$$

$$\begin{aligned} \tilde{A}^{(j)} &= T^{-1} A T \\ \tilde{u} &= T' u \\ \tilde{v} &= T^{-1} v \end{aligned}$$

Quasi-HMM realization (correctness of the algorithm)

Proposition 1. Let $\theta^o = (d, k, u, v, A^{(j)} : j \in [d])$ be a minimal order quasi realization for the output process. If the matrix $H^{(0)}$ has rank k , then the Algorithm returns an operator model $\tilde{\theta}^o$ that is equivalent to θ^o up to linear transformation.

Assume the output process can be realized by an HMM model of order k^h . The order k^h is always lower bounded by the order of the minimal order quasi realization k .

- ♦ Sketch of the proof

Unique up to linear transformation

$$H^{(0)} = UV', \quad U, V \text{ both of full rank } k$$

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix} \quad T = \begin{bmatrix} \tilde{u}' \tilde{A}^{(1)} \dots \tilde{A}^{(1)} \\ \tilde{u}' \tilde{A}^{(1)} \dots \tilde{A}^{(2)} \\ \vdots \\ \tilde{u}' \tilde{A}^{(d)} \dots \tilde{A}^{(d)} \end{bmatrix}$$

$$\tilde{A}^{(j)} = T^{-1} A T$$

$$\tilde{u} = T'u$$

$$\tilde{v} = T^{-1}v$$

$$V' = T^{-1} [A^{(1)} \dots A^{(1)} v, \dots, A^{(d)} \dots A^{(d)} v] = [\tilde{A}^{(1)} \dots \tilde{A}^{(1)} v, \dots, \tilde{A}^{(d)} \dots \tilde{A}^{(d)} \tilde{v}]$$

Quasi-HMM realization (correctness of the algorithm)

Proposition 1. Let $\theta^o = (d, k, u, v, A^{(j)} : j \in [d])$ be a minimal order quasi realization for the output process. If the matrix $H^{(0)}$ has rank k , then the Algorithm returns an operator model $\tilde{\theta}^o$ that is equivalent to θ^o up to linear transformation.

Assume the output process can be realized by an HMM model of order k^h . The order k^h is always lower bounded by the order of the minimal order quasi realization k .

- ♦ Sketch of the proof

$$H^{(0)} = UV', \quad U, V \text{ both of full rank } k$$

$$U = \begin{bmatrix} u' A^{(1)} \dots A^{(1)} \\ u' A^{(1)} \dots A^{(2)} \\ \vdots \\ u' A^{(d)} \dots A^{(d)} \end{bmatrix}$$

$T = \begin{bmatrix} T^{(1)} & \cdots & T^{(d)} \end{bmatrix}$

Computation complexity
 $\mathcal{O}(d(d^n)^3)$

Want poly time algorithm
 $n \sim \mathcal{O}(\log_d k)$

$$V' = T^{-1} [A^{(1)} \dots A^{(1)} v, \dots, A^{(d)} \dots A^{(d)} v] = [\tilde{A}^{(1)} \dots \tilde{A}^{(1)} \tilde{v}, \dots, \tilde{A}^{(d)} \dots \tilde{A}^{(d)} \tilde{v}]$$

$$\begin{aligned} \tilde{A}^{(j)} &= T^{-1} A T \\ \tilde{u} &= T' u \\ \tilde{v} &= T^{-1} v \end{aligned}$$

Algorithm 2 (correctness)

Theorem 3: (1) $M = A \otimes B \otimes C$

If both A and B have full column rank, then with probability one (over random projections), Algorithm 1 uniquely recovers the factors A, B, C up to state permutation. Moreover, Theorem 2.1 can be applied to recover (Q, O) .

Algorithm 2 (correctness)

Theorem 3: (1) $M = A \otimes B \otimes C$

If both A and B have full column rank, then with probability one (over random projections), Algorithm 1 uniquely recovers the factors A, B, C up to state permutation. Moreover, Theorem 2.1 can be applied to recover (Q, O) .

$$A = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O Q) \dots)Q)Q)Q}_{n \text{ } n},$$
$$B = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O \tilde{Q}) \dots)\tilde{Q})\tilde{Q})\tilde{Q}}_{n \text{ } n},$$
$$C = O \text{diag}(\pi)$$

Algorithm 2 (correctness)

Theorem 3: (1) $M = A \otimes B \otimes C$

If both A and B have full column rank, then with probability one (over random projections), Algorithm 1 uniquely recovers the factors A, B, C up to state permutation. Moreover, Theorem 2.1 can be applied to recover (Q, O) .

$$A = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O Q) \dots)Q)Q)Q}_{n \text{ } n},$$
$$B = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O \tilde{Q}) \dots)\tilde{Q})\tilde{Q})\tilde{Q}}_{n \text{ } n},$$
$$C = O \text{diag}(\pi)$$

Q needs to be of full rank k

Increasing n can boost the rank of A,B

How large N=2n+1 needs to be?

Algorithm 2 (correctness)

Theorem 3: (1) $M = A \otimes B \otimes C$

If both A and B have full column rank, then with probability one (over random projections), Algorithm 1 uniquely recovers the factors A, B, C up to state permutation. Moreover, Theorem 2.1 can be applied to recover (Q, O) .

$$\Theta_{(d,k)}^h = \{\theta^h : \text{order } k, \text{ alphabet size } d\}$$

There exists a measure zero set $\mathcal{E} \in \Theta_{(d,k)}^h$ such that for all HMMs in the set $\Theta_{(d,k)}^h \setminus \mathcal{E}$, for $N = 2n + 1$, with

$$n > 2 \log_d(k),$$

the matrices A and B have full column rank, thus $\mathcal{S}^{(N)}$ is sufficient for learning the minimal order HMM of the output process.

Algorithm 2 (correctness)

Theorem 3: (1) $M = A \otimes B \otimes C$

If both A and B have full column rank, then with probability one (over random projections), Algorithm 1 uniquely recovers the factors A, B, C up to state permutation. Moreover, Theorem 2.1 can be applied to recover (Q, O) .

$\Theta_{(d,k)}^h = \{\theta^h : \text{order } k, \text{ alphabet}$

Polynomial computational complexity

$$\mathcal{O}(d(d^n)^3) \sim \mathcal{O}(dk^6)$$

There exists a measure zero set $\mathcal{E} \in \Theta_{(d,k)}^h$ such that for all HMMs in the set $\Theta_{(d,k)}^h \setminus \mathcal{E}$, for $N = 2n + 1$, with

$$n > 2 \log_d(k),$$

the matrices A and B have full column rank, thus $\mathcal{S}^{(N)}$ is sufficient for learning the minimal order HMM of the output process.

Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) \quad M \Leftrightarrow A, B, C \rightarrow (2) \quad A, B, C \Leftrightarrow (Q, O)$$

Proposition 2: sufficient conditions for invertability of (1)

Consider the matrices A, B, C parameterized by a minimal realization (Q, O) . If the following deterministic conditions on the Kruskal rank are satisfied:

$$\text{krank}(A) + \text{krank}(B) + \text{krank}(C) \geq 2k + 2, \quad (1)$$

then M can be uniquely (up to common column permutation) decomposed into A, B, C , and k can be determined thereby.

Minimal HMM realization (identifiability)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) \quad M \Leftrightarrow A, B, C \rightarrow (2) \quad A, B, C \Leftrightarrow (Q, O)$$

Proposition 2: sufficient conditions for invertability of (1)

Consider the matrices A, B, C parameterized by a minimal realization (Q, O) . If the following deterministic conditions on the Kruskal rank are satisfied:

$$\text{krank}(A) + \text{krank}(B) + \text{krank}(C) \geq 2k + 2, \quad (1)$$

then M can be uniquely (up to common column permutation) decomposed into A, B, C , and k can be determined thereby.

$$A = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O Q) \dots) Q) Q) Q}_{n \text{ blocks}},$$

$$B = \underbrace{(O \odot (O \odot (O \odot \dots (O \odot O \tilde{Q}) \dots) \tilde{Q}) \tilde{Q}) \tilde{Q}}_{n \text{ blocks}},$$

$$C = O \text{diag}(\pi)$$

Algorithm 4 (random instance check)

$\Theta_{(d,k,r)}^h = \{\theta^h : \text{order } k, \text{ alphabet size } d, Q \text{ is of rank } r < k\}$

1. Randomly pick an HMM model $(d, k, Q, O) \in \Theta_{(d,k,r)}^h$.
2. Construct matrices $A = OQ$, $B = O\tilde{Q}$, $C = ODiag(\pi)$,
where π is the stationary distribution, and $\tilde{Q} = Diag(\pi)Q'Diag(\pi)^{-1}$.
3. Run Algorithm 2.a with input $M = A \otimes B \otimes C$.
4. Return “yes” if the algorithm recovers A, B, C up to state permutation,
and “no” otherwise.

Algorithm 4 (random instance check)

$$\Theta_{(d,k,r)}^h = \{\theta^h : \text{order } k, \text{ alphabet size } d, Q \text{ is of rank } r < k\}$$

1. Randomly pick an HMM model $(d, k, Q, O) \in \Theta_{(d,k,r)}^h$.
2. Construct matrices $A = OQ$, $B = O\tilde{Q}$, $C = ODiag(\pi)$,
where π is the stationary distribution, and $\tilde{Q} = Diag(\pi)Q'Diag(\pi)^{-1}$.
3. Run Algorithm 2.a with input $M = A \otimes B \otimes C$.
4. Return “yes” if the algorithm recovers A, B, C up to state permutation,
and “no” otherwise.

Theorem 4: Correctness of Algorithm 3

For given $d \geq k$ and $n = 1$, if Algorithm 2.b returns “yes”, then there exists a measure zero set $\mathcal{E} \in \Theta_{(d,r,k)}^h$ such that Algorithm 2.a can identify all minimal HMM realizations in the set $\Theta_{(d,r,k)}^h \setminus \mathcal{E}$.

Moreover, if the latter is true, Algorithm 2.b returns “yes” with probability
1. Theorem 1 (2) can be applied to recover (Q, O) for this class of HMMs.

Minimal HMM realization (summary)

$$N \rightarrow \mathcal{S}^{(N)} \rightarrow (1) M \Leftrightarrow A, B, C \rightarrow (2) A, B, C \Leftrightarrow (Q, O)$$

- ♦ **Identifiability:** if (1) (2) are bijective mappings
 - ✓ How large N is so that the minimal HMM is identifiable from $\mathcal{S}^{(N)}$?
- ♦ **Algorithms:** efficiently invert (1) (2)
 - ✓ For (1) : Algorithm 2 (Simultaneous diagonalization)
 - If Q is full rank
 - Algorithm 3 (Foobi)
 - If O has full column rank, and check a random instance
 - ✓ For (2) : Theorem 2 (Linear inversion)
 - If Q or O has full column rank