

---

# **PATTERN RECOGNITION AND MACHINE LEARNING**

## **CHAPTER 13: SEQUENTIAL DATA**

---

# Learning Objectives

- 1、 What are hidden Markov models (HMMs)?
- 2、 What is the EM scheme for HMMs?
- 3、 What are Forward-Backward and Sum-Product Algorithms?
- 4、 What are Viterbi and Max-Product Algorithms?
- 5、 What are linear dynamic systems?
- 6、 What are Kalman and particle filters?
- 7、 How to learn linear dynamic system models?
- 8、 What are RNN and LSTM?

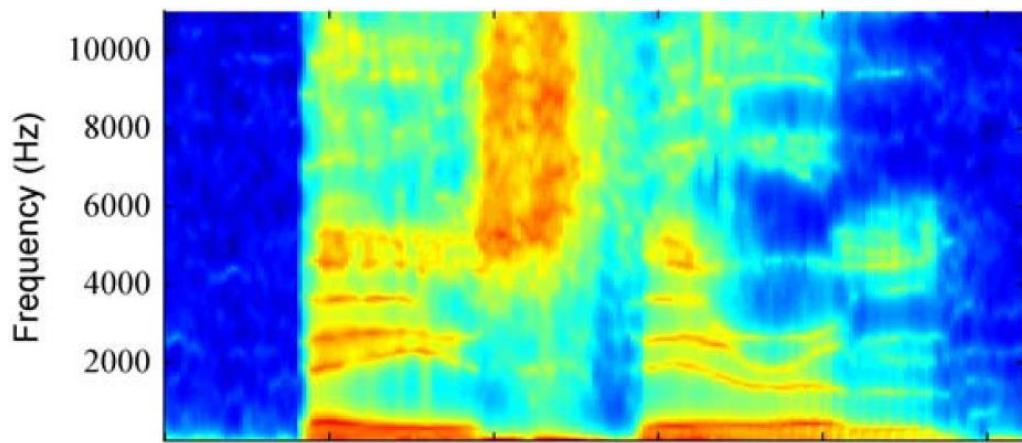
# Outlines

---

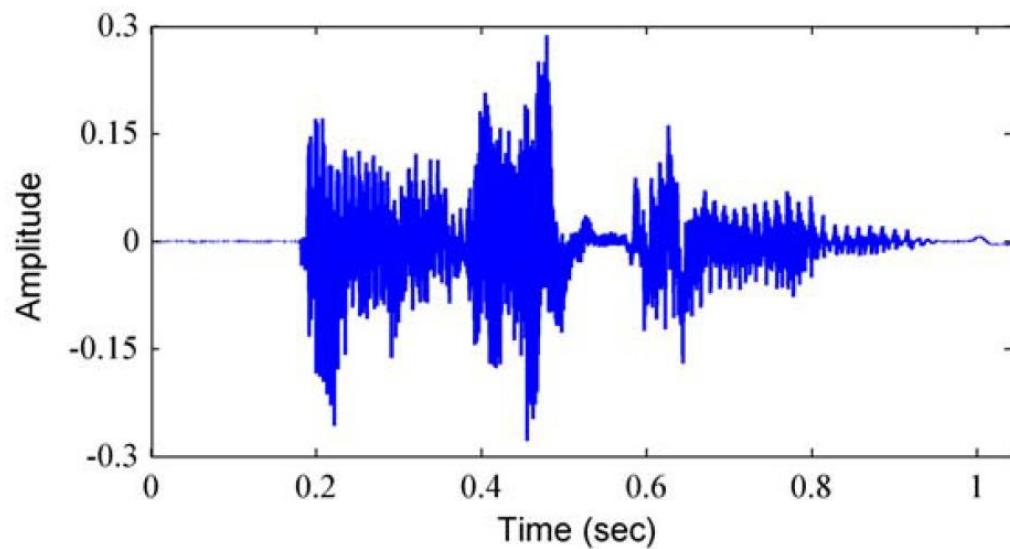
- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-

# Sequential Data

---

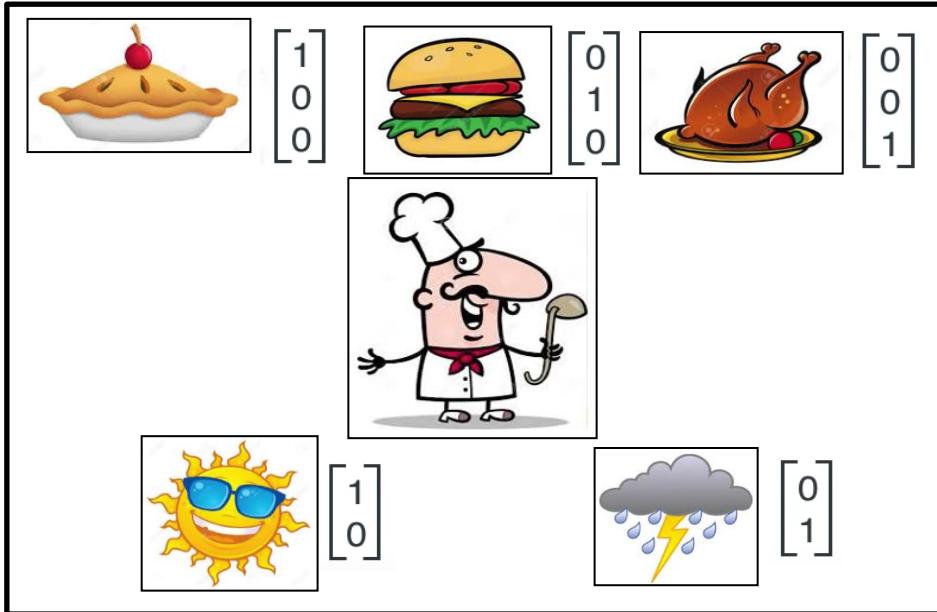


| b |   ey |   z |   th |   ih |   er |   em |  
| Bayes' |                  | Theorem |



Spectral and temporal signals  
of Bayes' Theorem

# Conditional Data

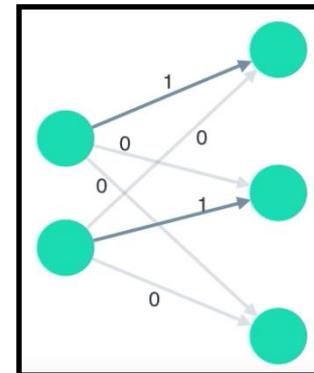


$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

A diagram showing a sun wearing sunglasses next to a pie. This represents the result of the matrix multiplication where the condition is 1 (sunny).

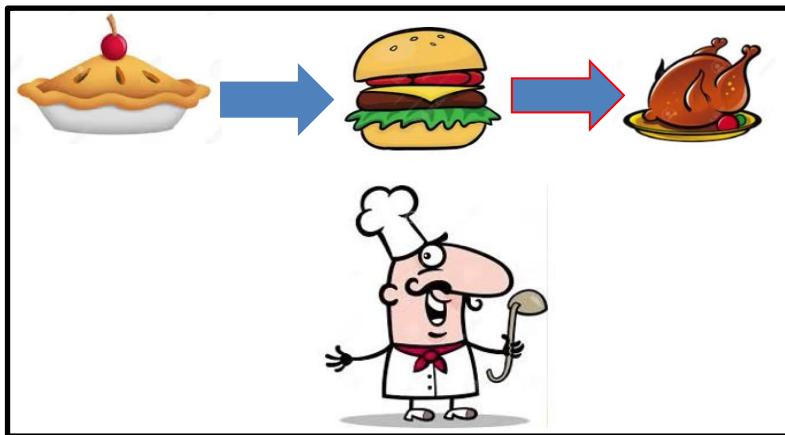
$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

A diagram showing a sun wearing sunglasses next to a burger. This represents the result of the matrix multiplication where the condition is 0 (rainy).



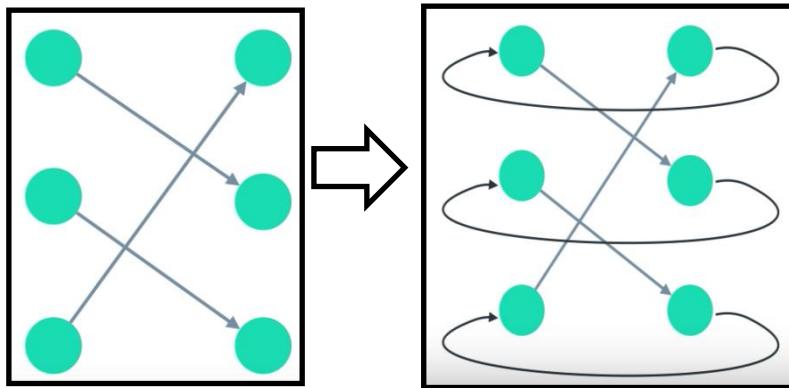
# Sequential Data

---



$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{array}{c} \text{Pie} \\ \text{Burger} \\ \text{Roasted Chicken} \end{array} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{array}{c} \text{Pie} \\ \text{Burger} \\ \text{Roasted Chicken} \end{array}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{array}{c} \text{Burger} \\ \text{Roasted Chicken} \end{array} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{array}{c} \text{Burger} \\ \text{Roasted Chicken} \end{array}$$



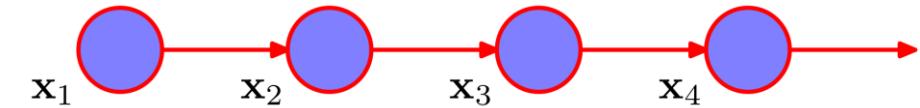
$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \begin{array}{c} \text{Roasted Chicken} \\ \text{Pie} \end{array} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{array}{c} \text{Roasted Chicken} \\ \text{Pie} \end{array}$$

# Markov Models

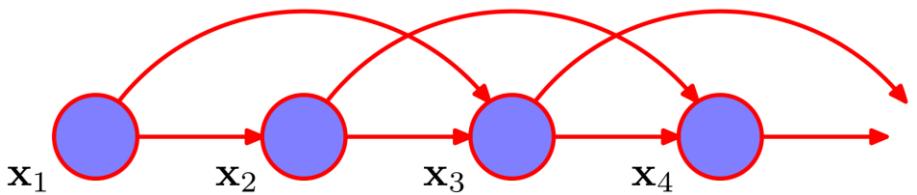
---

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) \prod_{n=2}^N p(\mathbf{x}_n | \mathbf{x}_{n-1})$$

$$p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) = p(\mathbf{x}_n | \mathbf{x}_{n-1})$$



A second-order Markov chain



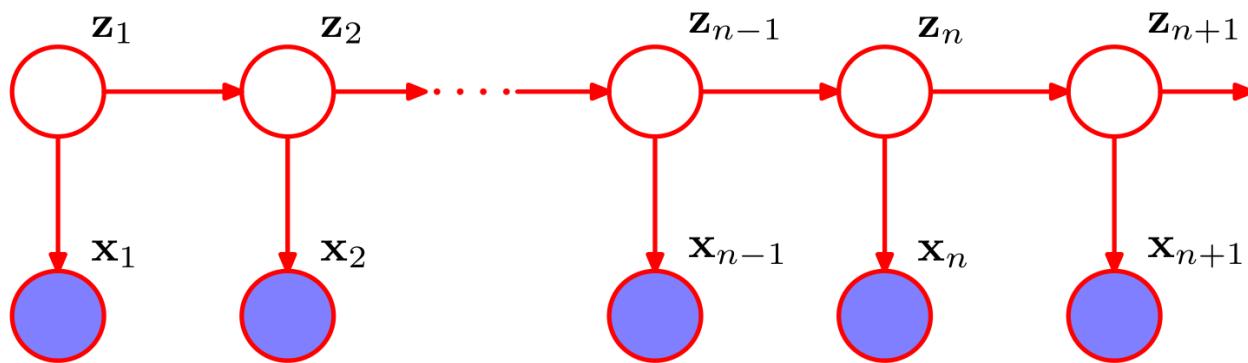
$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1)p(\mathbf{x}_2 | \mathbf{x}_1) \prod_{n=3}^N p(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{x}_{n-2})$$

# Hidden Markov Models

---

Using a Markov chain of latent variables

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{z}_1, \dots, \mathbf{z}_N) = p(\mathbf{z}_1) \left[ \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}) \right] \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n)$$



For continuous variables, we can use linear-Gaussian conditional distributions

---

# Latent Markov Model

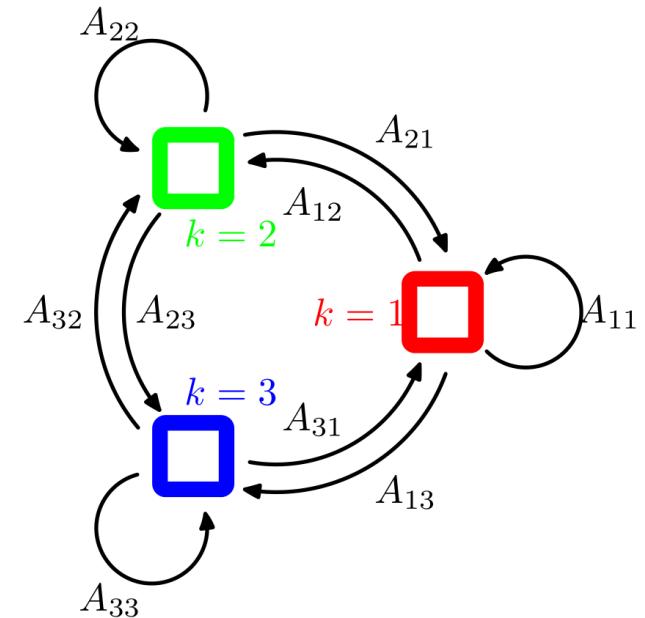
Conditional distribution for latent variable

$$p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$$

$$p(\mathbf{z}_1 | \boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{z_{1k}}$$

$$\sum_k \pi_k = 1$$

A means transition probabilities



As in the case of a standard mixture model, the latent variables are the discrete multinomial variables  $\mathbf{z}_n$  using the 1-of-K coding scheme

A model whose latent variables have three possible states corresponding to the three boxes. The black lines denote the elements of the transition matrix  $A_{jk}$

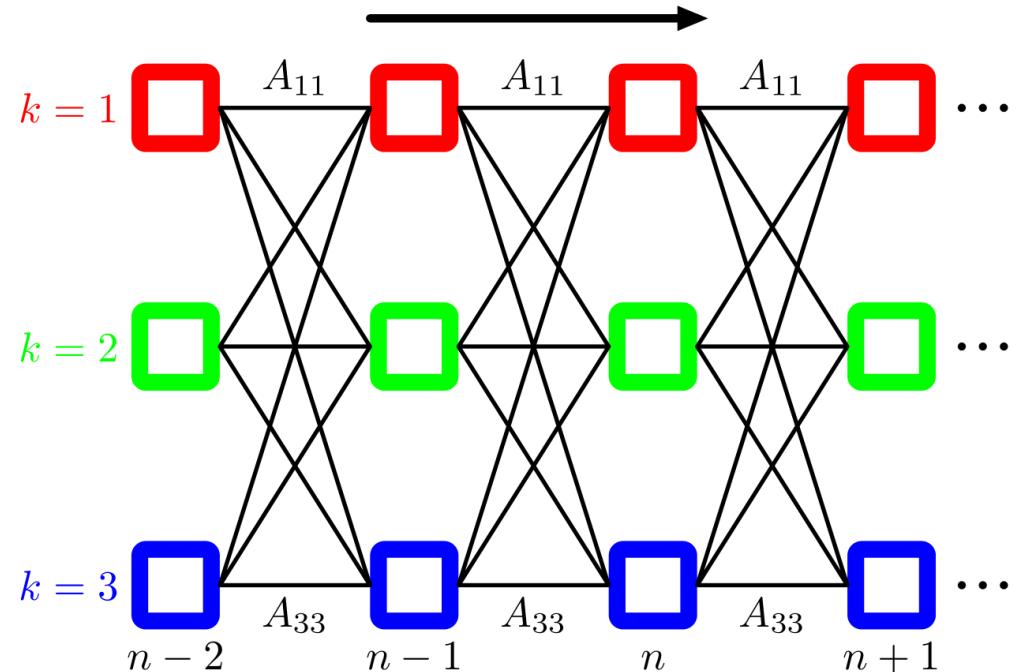
# Latent State Lattice

---

Latent states lattice:  
representing the transitions  
between latent states

Each column of this diagram  
corresponds to one of the  
latent variables  $\mathbf{z}_n$

Each row of this diagram  
corresponds to one state of  
the latent variables  $\mathbf{z}_n$



# Hidden Markov Models

---

Getting emission probabilities from latent variable

$$p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\phi}) = \prod_{k=1}^K p(\mathbf{x}_n | \boldsymbol{\phi}_k)^{z_{nk}} \quad p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}) = \prod_{k=1}^K \prod_{j=1}^K A_{jk}^{z_{n-1,j} z_{nk}}$$

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\theta}) = p(\mathbf{z}_1 | \boldsymbol{\pi}) \left[ \prod_{n=2}^N p(\mathbf{z}_n | \mathbf{z}_{n-1}, \mathbf{A}) \right] \prod_{m=1}^N p(\mathbf{x}_m | \mathbf{z}_m, \boldsymbol{\phi})$$

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

$$\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$$

$$\boldsymbol{\theta} = \{\boldsymbol{\pi}, \mathbf{A}, \boldsymbol{\phi}\}$$

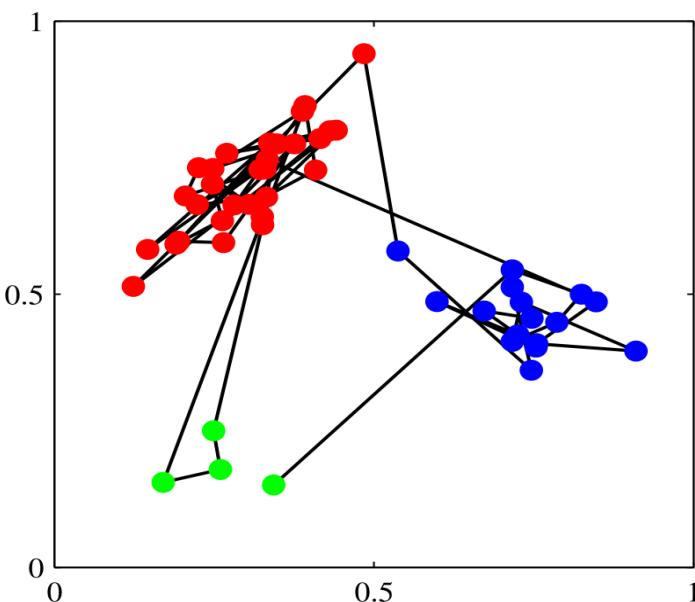
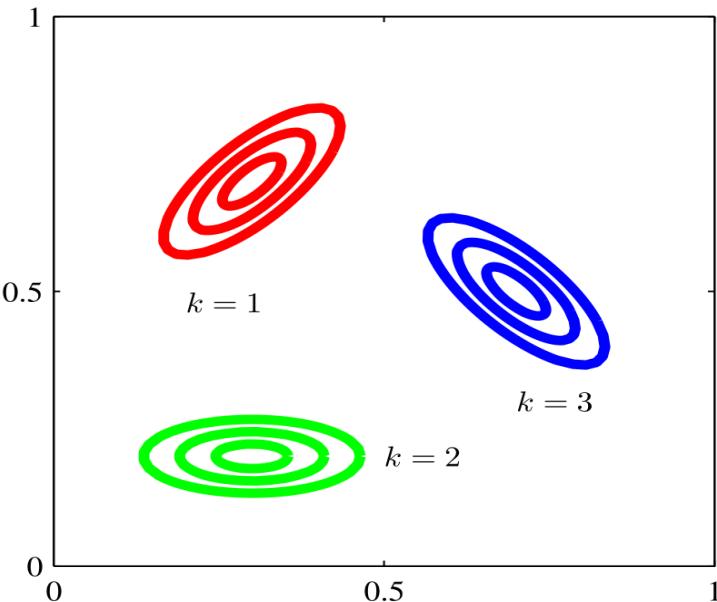
---

# Hidden Markov Model Data Sequence

A better understanding of the hidden Markov model

First choose the initial latent variable  $\mathbf{z}_1 \ \mathbf{x}_1 \ \pi_k$  state  $j$

Then, choose the state of the variable  $\mathbf{z}_2 \ p(\mathbf{z}_2|\mathbf{z}_1)$  state  $k \ A_{jk} \ k = 1, \dots, K$

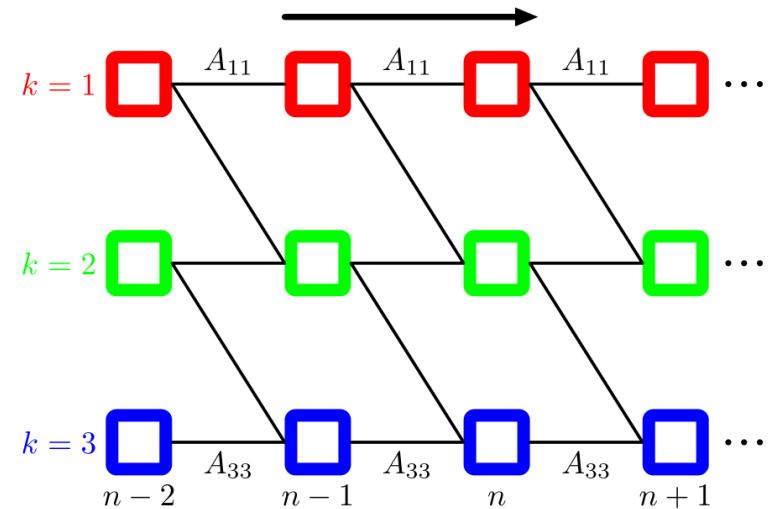
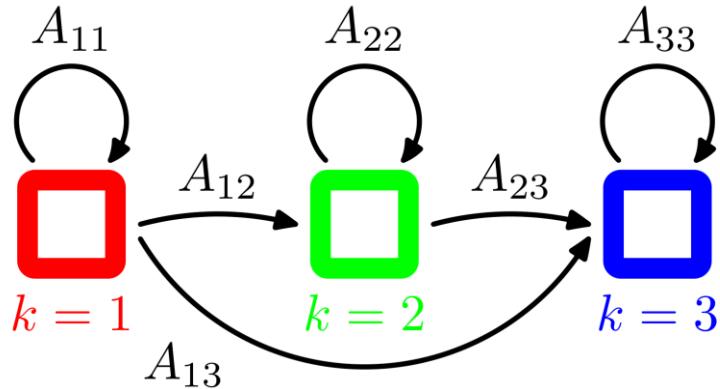


The number of latent variable states is 3;  
the dimension of emission variable is 2

A sequence of 50 samples

# Hidden Markov Model Examples

One Example of variants of the standard HMM model

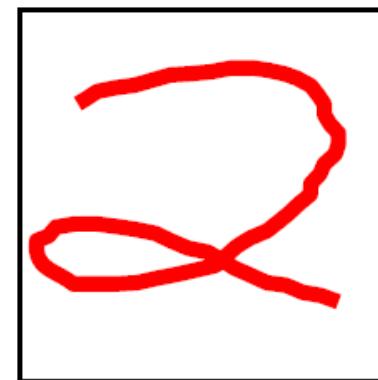
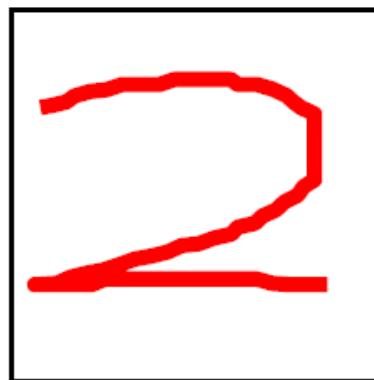


3-state left-to-right hidden Markov model

Lattice diagram for a 3-state left- to-right HMM

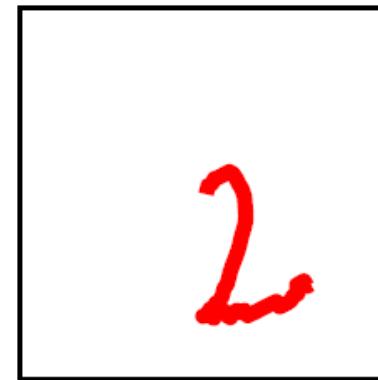
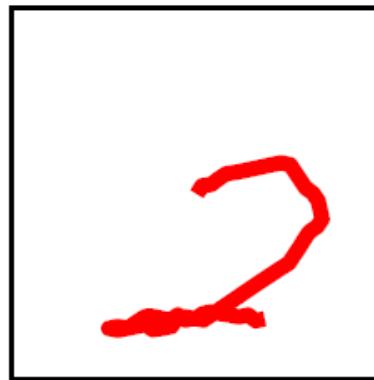
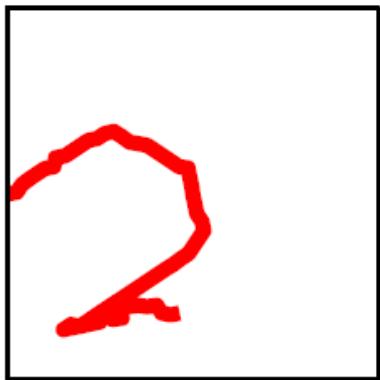
# HMM Applications

---



Handwritten digits

Synthetic digits using HMMs which are trained by using handwritten digits



# Outlines

---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi Algorithm
  - Linear Dynamics Systems
  - Kalman and Particle Filters
  - RNN and LSTM
-

# Three Problems for HMMs

---

## □ Evaluation:

Given a HMM model  $\theta = \{\pi, \mathbf{A}, \phi\}$ , what is likelihood of an observation sequence  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  generated by that model?

## □ Learning:

What is the most likely HMM model for an observation sequence  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ ?

## □ Decoding:

Given a HMM model  $\theta = \{\pi, \mathbf{A}, \phi\}$ , what is the most likely latent sequence  $\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$  for an observation sequence  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  ?

---

# Expectation of Latent Variables

---

Expectation of latent variables

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$\mathbf{X}$ : all data samples

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{nk}$$

$$\xi(z_{n-1,j}, z_{nk}) = \mathbb{E}[z_{n-1,j} z_{nk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{n-1,j} z_{nk}$$

# Maximum Likelihood of HMMs

---

Maximum likelihood for the HMM

$$p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$

EM algorithm to find an efficient framework

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})$$



$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) &= \sum_{k=1}^K \gamma(z_{1k}) \ln \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1,j}, z_{nk}) \ln A_{jk} \\ &\quad + \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n | \boldsymbol{\phi}_k). \end{aligned}$$

---

# EM Learning of HMMs (I)

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \sum \gamma(\mathbf{z}) z_{nk}$$

$$\xi(z_{n-1,j}, z_{nk}) = \mathbb{E}[z_{n-1,j} z_{nk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{n-1,j} z_{nk}$$

$$p(\mathbf{x} | \mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)^{z_k}$$

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})}$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})}$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T}{\sum_{n=1}^N \gamma(z_{nk})}$$

# EM Learning of HMMs (II)

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$$

$$\gamma(z_{nk}) = \mathbb{E}[z_{nk}] = \sum \gamma(\mathbf{z}) z_{nk}$$

$$\xi(z_{n-1,j}, z_{nk}) = \mathbb{E}[z_{n-1,j} z_{nk}] = \sum_{\mathbf{z}} \gamma(\mathbf{z}) z_{n-1,j} z_{nk}$$

$$p(\mathbf{x} | \mathbf{z}) = \prod_{i=1}^D \prod_{k=1}^K \mu_{ik}^{x_i z_k}$$

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})}$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=1}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})}$$

$$\mu_{ik} = \frac{\sum_{n=1}^N \gamma(z_{nk}) x_{ni}}{\sum_{n=1}^N \gamma(z_{nk})}$$

# Outlines

---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-

# Forward Recursion (I)

## □ Forward recursion

$$\begin{aligned} p(\mathbf{X}|\mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) \\ &\quad p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \end{aligned}$$

$$p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{x}_n, \mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_n)$$

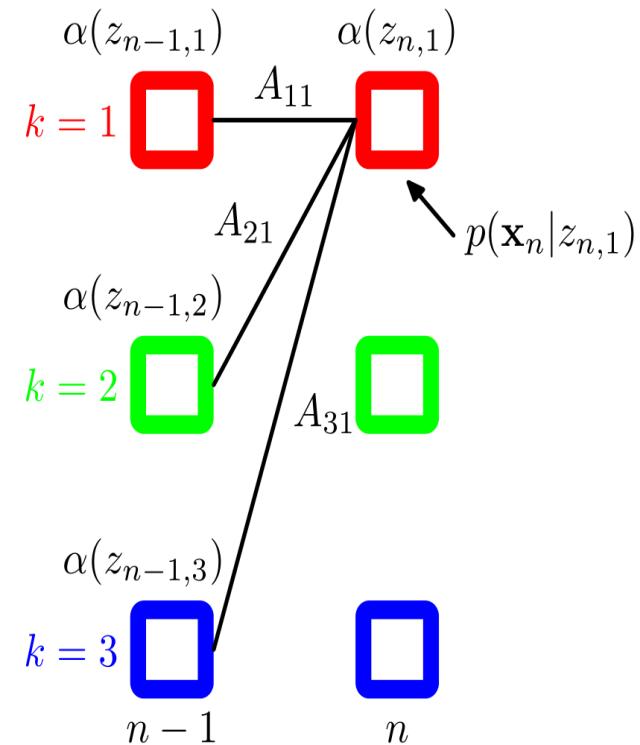
$$p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1})$$

$$p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n, \mathbf{z}_{n+1}) = p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1})$$

$$\begin{aligned} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}, \mathbf{x}_{n+1}) &= p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) \\ p(\mathbf{X} | \mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}) \\ &\quad p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \end{aligned}$$

$$p(\mathbf{x}_{N+1} | \mathbf{X}, \mathbf{z}_{N+1}) = p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1})$$

$$p(\mathbf{z}_{N+1} | \mathbf{z}_N, \mathbf{X}) = p(\mathbf{z}_{N+1} | \mathbf{z}_N)$$



$$\alpha(\mathbf{z}_1) = p(\mathbf{x}_1, \mathbf{z}_1) = p(\mathbf{z}_1)p(\mathbf{x}_1 | \mathbf{z}_1) = \prod_{k=1}^K \{\pi_k p(\mathbf{x}_1 | \phi_k)\}^{z_{1k}}$$

# Forward Recursion (II)

---

$$\begin{aligned}\alpha(\mathbf{z}_n) &= p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) \\&= p(\mathbf{x}_1, \dots, \mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) \\&= p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_n) p(\mathbf{z}_n) \\&= p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n) \\&= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}, \mathbf{z}_n) \\&= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1}) \\&= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1} | \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1}) p(\mathbf{z}_{n-1}) \\&= p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} p(\mathbf{x}_1, \dots, \mathbf{x}_{n-1}, \mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})\end{aligned}$$

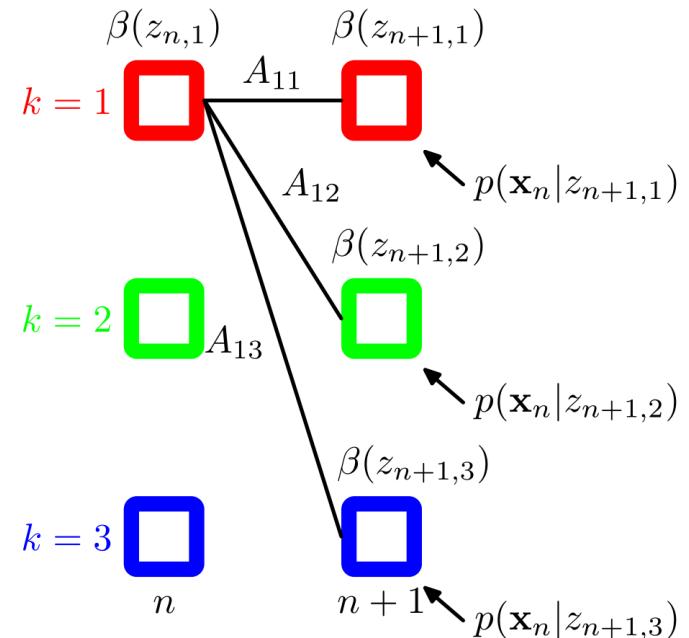
$$\boxed{\alpha(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \alpha(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})}$$

---

# Backward Recursion

## □ Backward recursion

$$\begin{aligned}\beta(\mathbf{z}_n) &= p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N, \mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n, \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n) \\ &= \sum_{\mathbf{z}_{n+1}} p(\mathbf{x}_{n+2}, \dots, \mathbf{x}_N | \mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)\end{aligned}$$



$$\boxed{\beta(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \beta(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)}$$

# Forward-Backward Estimation

The method of evaluating the quantities of  $\gamma(z_{nk})$   $\xi(z_{n-1,j}, z_{nk})$

$$\begin{aligned}\gamma(\mathbf{z}_n) &= p(\mathbf{z}_n | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{z}_n)p(\mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{X})} = \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})}\end{aligned}$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X})$$

$$= \frac{\alpha(\mathbf{z}_{n-1})p(\mathbf{x}_n | \mathbf{z}_n)p(\mathbf{z}_n | \mathbf{z}_{n-1})\beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

$$\begin{aligned}\alpha(\mathbf{z}_n) &\equiv p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n) \\ \beta(\mathbf{z}_n) &\equiv p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)\end{aligned}$$

$$p(\mathbf{X}) = \sum_{\mathbf{z}_n} \alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)$$

$$p(\mathbf{X}) = \sum_{\mathbf{z}_N} \alpha(\mathbf{z}_N)$$

Evaluation

# Observation Prediction

---

$$\begin{aligned} p(\mathbf{x}_{N+1} | \mathbf{X}) &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1}, \mathbf{z}_{N+1} | \mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) p(\mathbf{z}_{N+1} | \mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1}, \mathbf{z}_N | \mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1} | \mathbf{z}_N) p(\mathbf{z}_N | \mathbf{X}) \\ &= \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1} | \mathbf{z}_N) \frac{p(\mathbf{z}_N, \mathbf{X})}{p(\mathbf{X})} \\ &= \frac{1}{p(\mathbf{X})} \sum_{\mathbf{z}_{N+1}} p(\mathbf{x}_{N+1} | \mathbf{z}_{N+1}) \sum_{\mathbf{z}_N} p(\mathbf{z}_{N+1} | \mathbf{z}_N) \alpha(\mathbf{z}_N) \end{aligned}$$

---

# Sum-Product v.s. Max-Product

---

## □ Sum-Product Algorithm (evaluation)

- ✓ Compute the joint distribution from the Product
- ✓ Infer marginal distributions from the Sum

$$p(x_1, x_2) = \sum_{x_3} p(x_1, x_3)p(x_2, x_3)$$

## □ Max-Product Algorithm (decoding)

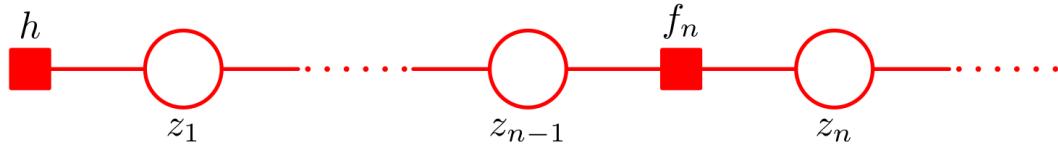
- ✓ Compute the joint distribution from the Product
- ✓ Perform ML estimation from the Max

$$x_1^* = \max_{x_1} p(x_1, x_3)p(x_2, x_1)$$

---

# Sum-Product for HMMs

- Transforming the directed graph into a factor graph



$$h(\mathbf{z}_1) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1)$$

$$f_n(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)$$

$$\mu_{\mathbf{z}_{n-1} \rightarrow f_n}(\mathbf{z}_{n-1}) = \mu_{f_{n-1} \rightarrow \mathbf{z}_{n-1}}(\mathbf{z}_{n-1})$$

$$\mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) = \sum_{\mathbf{z}_{n-1}} f_n(\mathbf{z}_{n-1}, \mathbf{z}_n) \mu_{\mathbf{z}_{n-1} \rightarrow f_n}(\mathbf{z}_{n-1})$$

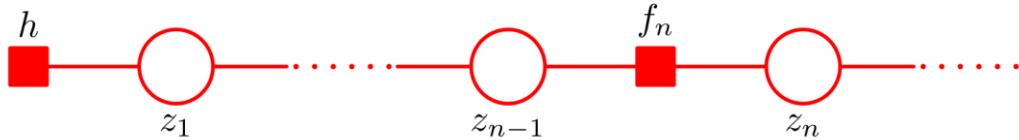
$$\boxed{\mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) = \sum_{\mathbf{z}_{n-1}} f_n(\mathbf{z}_{n-1}, \mathbf{z}_n) \mu_{f_{n-1} \rightarrow \mathbf{z}_{n-1}}(\mathbf{z}_{n-1})}$$

$$\boxed{\alpha(\mathbf{z}_n) = \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n)}$$

# Sum-Product for HMMs

---

- Transforming the directed graph into a factor graph



$$\mu_{f_{n+1} \rightarrow f_n}(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} f_{n+1}(\mathbf{z}_n, \mathbf{z}_{n+1}) \mu_{f_{n+2} \rightarrow f_{n+1}}(\mathbf{z}_{n+1})$$

$$\beta(\mathbf{z}_n) = \mu_{f_{n+1} \rightarrow \mathbf{z}_n}(\mathbf{z}_n)$$

$$p(\mathbf{z}_n, \mathbf{X}) = \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n) \mu_{f_{n+1} \rightarrow \mathbf{z}_n}(\mathbf{z}_n) = \alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)$$

$$\gamma(\mathbf{z}_n) = \frac{p(\mathbf{z}_n, \mathbf{X})}{p(\mathbf{X})} = \frac{\alpha(\mathbf{z}_n) \beta(\mathbf{z}_n)}{p(\mathbf{X})}$$

# Scaling of HMMs

## □ Scaling factors for computational practice

$$\alpha(\mathbf{z}_n) = p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_n)$$

$$\widehat{\alpha}(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{\alpha(\mathbf{z}_n)}{p(\mathbf{x}_1, \dots, \mathbf{x}_n)}$$

scaling factors defined by conditional distributions over the observed variables

$$c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

$$c_n \widehat{\alpha}(\mathbf{z}_n) = p(\mathbf{x}_n | \mathbf{z}_n) \sum_{\mathbf{z}_{n-1}} \widehat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{z}_n | \mathbf{z}_{n-1})$$

similarly define re-scaled variables

$$\widehat{\beta}(\mathbf{z}_n) = \frac{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{z}_n)}{p(\mathbf{x}_{n+1}, \dots, \mathbf{x}_N | \mathbf{x}_1, \dots, \mathbf{x}_n)}$$

$$c_{n+1} \widehat{\beta}(\mathbf{z}_n) = \sum_{\mathbf{z}_{n+1}} \widehat{\beta}(\mathbf{z}_{n+1}) p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) p(\mathbf{z}_{n+1} | \mathbf{z}_n)$$



$$\gamma(\mathbf{z}_n) = \widehat{\alpha}(\mathbf{z}_n) \widehat{\beta}(\mathbf{z}_n)$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = c_n^{-1} \widehat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{n-1}) \widehat{\beta}(\mathbf{z}_n)$$

# Outlines

---

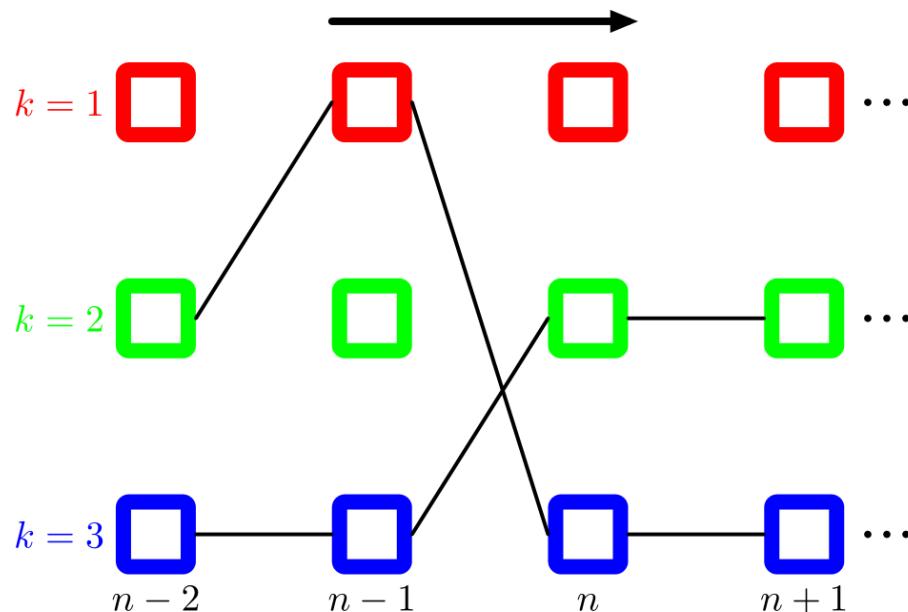
- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-

# Latent Sequence Estimation

---

## □ Decoding:

Given a HMM model  $\theta = \{\pi, \mathbf{A}, \phi\}$ , what is the most likely latent sequence  $\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$  for an observation sequence  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  ?



# Viterbi Algorithm

---

$$\omega(\mathbf{z}_n) = \max_{\mathbf{z}_1, \dots, \mathbf{z}_{n-1}} \ln p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{z}_1, \dots, \mathbf{z}_n)$$

$$\omega(\mathbf{z}_{n+1}) = \max_{\mathbf{z}_1, \dots, \mathbf{z}_n} \ln p(\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{x}_{n+1}, \mathbf{z}_1, \dots, \mathbf{z}_n, \mathbf{z}_{n+1})$$

$$\omega(\mathbf{z}_{n+1}) = \ln p(\mathbf{x}_{n+1} | \mathbf{z}_{n+1}) + \max_{\mathbf{z}_n} \{\ln p(\mathbf{z}_{n+1} | \mathbf{z}_n) + \omega(\mathbf{z}_n)\}$$

$$\omega(\mathbf{z}_1) = \ln p(\mathbf{z}_1) + \ln p(\mathbf{x}_1 | \mathbf{z}_1)$$

---

# Viterbi Algorithm

---

- Note that maximization over  $\mathbf{z}_n$  must be performed for each of K possible values of  $\mathbf{z}_{n+1}$
- Denote this function by  $\psi(k_n)$ , where  $k \in \{1, \dots, K\}$
- Once we find the most probable value of  $\mathbf{z}_N$ , we can trackback along the chain

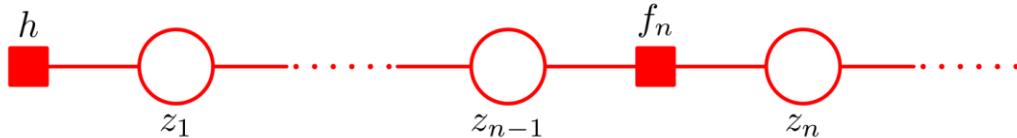
$$k_n^{\max} = \psi(k_{n+1}^{\max})$$

- Reduce the computational cost from  $O(K^N)$  to  $O(KN)$
-

# Max-Product for HMMs

---

- Transforming the directed graph into a factor graph



$$h(\mathbf{z}_1) = p(\mathbf{z}_1)p(\mathbf{x}_1|\mathbf{z}_1)$$

$$f_n(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_n|\mathbf{z}_{n-1})p(\mathbf{x}_n|\mathbf{z}_n)$$

$$\mu_{\mathbf{z}_n \rightarrow f_{n+1}}(\mathbf{z}_n) = \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n)$$

$$\mu_{f_{n+1} \rightarrow \mathbf{z}_{n+1}}(\mathbf{z}_{n+1}) = \max_{\mathbf{z}_n} \left\{ \ln f_{n+1}(\mathbf{z}_n, \mathbf{z}_{n+1}) + \mu_{\mathbf{z}_n \rightarrow f_{n+1}}(\mathbf{z}_n) \right\}$$

- As such, the  $f \rightarrow \mathbf{z}$  message will recursively be

$$\omega(\mathbf{z}_{n+1}) = \ln p(\mathbf{x}_{n+1}|\mathbf{z}_{n+1}) + \max_{\mathbf{z}_n} \left\{ \ln p(\mathbf{z}_{n+1}|\mathbf{z}_n) + \omega(\mathbf{z}_n) \right\}$$

$$\omega(\mathbf{z}_n) \equiv \mu_{f_n \rightarrow \mathbf{z}_n}(\mathbf{z}_n)$$

---

# Discriminative HMMs

---

## □ Using discriminative rather than Maximum Likelihood techniques

optimize the cross-entropy of  $R$  observation sequences,  $\mathbf{X}_r$ , and labels,  $m_r$

$$\sum_{r=1}^R \ln p(m_r | \mathbf{X}_r) \iff \sum_{r=1}^R \ln \left\{ \frac{p(\mathbf{X}_r | \boldsymbol{\theta}_r) p(m_r)}{\sum_{l=1}^M p(\mathbf{X}_r | \boldsymbol{\theta}_l) p(l_r)} \right\} \quad r = 1, \dots, R \\ m = 1, \dots, M$$

for each class there is a HMM,  $\boldsymbol{\theta}_m$

## □ Weakness of the first-order HMM:

- ✓ distribution of times for which the system remains in a given state
- ✓ poor at capturing long-range correlations between the observed variables

# Example

---

- Given an HMM and an observation sequence, how to perform evaluation and decoding

Transition A	Emission B	Hidden States Z	Observations X
$\begin{bmatrix} 0.6 & 0.3 \\ 0.4 & 0.7 \end{bmatrix}$	$\begin{bmatrix} 0.8 & 0.1 \\ 0.2 & 0.9 \end{bmatrix}$	{bull, bear}	{up, down}

If Z is stationary, then  $\pi = [3/7, 4/7]$ . We also can assume  $\pi = [1/2, 1/2]$

An observation sequence: {up, up, down}

# Evaluation (Sum-Product)

---

$$\alpha(z_1) = p(z_1, x_1) = p(x_1|z_1)p(z_1)$$

$$\boxed{x_1=\text{up}, z_1=\text{bull or bear}} \quad = \begin{bmatrix} 0.8 \\ 0.1 \end{bmatrix} \cdot \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.8 \times 0.5 \\ 0.1 \times 0.5 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0.05 \end{bmatrix}$$

$$\alpha(z_2) = p(z_2, x_1, x_2) = p(x_2|z_2) \sum_{z_1} p(z_2|z_1) \alpha(z_1)$$

$$\boxed{x_2=\text{up}, z_2=\text{bull or bear}} \quad = \begin{bmatrix} 0.8 \\ 0.1 \end{bmatrix} \cdot \begin{bmatrix} 0.6 \times 0.4 + 0.3 \times 0.05 \\ 0.4 \times 0.4 + 0.7 \times 0.05 \end{bmatrix} = \begin{bmatrix} 0.204 \\ 0.0195 \end{bmatrix}$$

$$\alpha(z_3) = p(z_3, x_1, x_2, x_3) = p(x_3|z_3) \sum_{z_2} p(z_3|z_2) \alpha(z_2)$$

$$\boxed{x_3=\text{down}, z_2=\text{bull or bear}} \quad = \begin{bmatrix} 0.2 \\ 0.9 \end{bmatrix} \cdot \begin{bmatrix} 0.6 \times 0.204 + 0.3 \times 0.0195 \\ 0.4 \times 0.204 + 0.7 \times 0.0195 \end{bmatrix} = \begin{bmatrix} 0.02565 \\ 0.085725 \end{bmatrix}$$

$$p(x_1, x_2, x_3) = \sum_{z_3} \alpha(z_3) = 0.111375$$

---

# Decoding (Max-Product)

$$\delta(z_1) = p(z_1, x_1) = p(x_1|z_1)p(z_1)$$

$$x_1=\text{up}, z_1=\text{bull or bear} \quad = \begin{bmatrix} 0.8 \\ 0.1 \end{bmatrix} \cdot \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.8 \times 0.5 \\ 0.1 \times 0.5 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0.05 \end{bmatrix}$$

$$\delta(z_2) = p(z_2, x_1, x_2) = p(x_2|z_2) \max_{z_1} p(z_2|z_1)\delta(z_1)$$

$$x_2=\text{up}, z_2=\text{bull or bear} \quad = \begin{bmatrix} 0.8 \\ 0.1 \end{bmatrix} \cdot \begin{bmatrix} 0.4 \times 0.6 \\ 0.4 \times 0.4 \end{bmatrix} = \begin{bmatrix} 0.192 \\ 0.016 \end{bmatrix} \leftarrow$$

$$\phi(z_2) = \arg \max_{z_1} p(z_2|z_1)\delta(z_1) = \begin{bmatrix} \text{bull} \rightarrow \text{bull} \\ \text{bull} \rightarrow \text{bear} \end{bmatrix} \leftarrow$$

$$\delta(z_3) = p(z_3, x_1, x_2, x_3) = p(x_3|z_3) \max_{z_2} p(z_3|z_2)\delta(z_2)$$

$$x_3=\text{down}, z_2=\text{bull or bear} \quad = \begin{bmatrix} 0.2 \\ 0.9 \end{bmatrix} \cdot \begin{bmatrix} 0.192 \times 0.6 \\ 0.192 \times 0.4 \end{bmatrix} = \begin{bmatrix} 0.02304 \\ 0.06912 \end{bmatrix} \leftarrow$$

$$\phi(z_3) = \arg \max_{z_2} p(z_3|z_2)\delta(z_2) = \begin{bmatrix} \text{bull} \rightarrow \text{bull} \\ \text{bull} \rightarrow \text{bear} \end{bmatrix} \leftarrow$$

Optimal solution:  
 $\text{bull} \rightarrow \text{bull} \rightarrow \text{bear}$

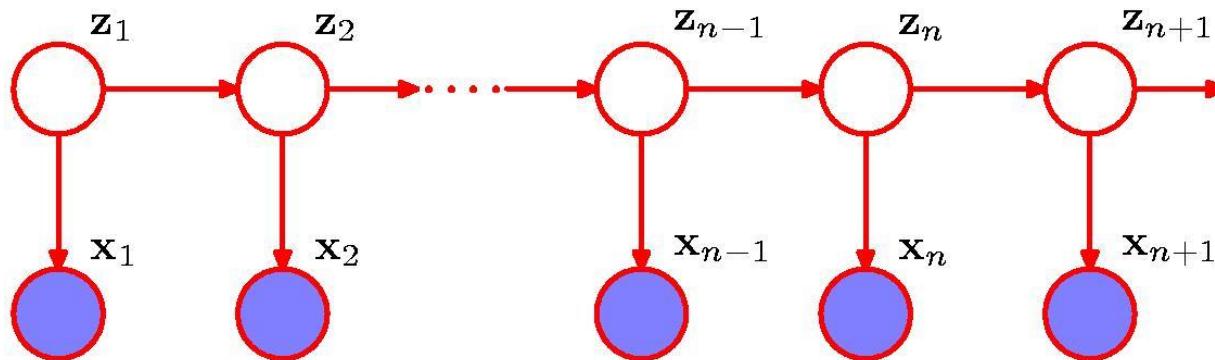
# Outlines

---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman and Particle Filters
  - RNN and LSTM
-

# Stochastic Linear Dynamical Systems

---



$$p(\mathbf{z}_n | \mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{z}_n | \mathbf{A}\mathbf{z}_{n-1}, \boldsymbol{\Gamma})$$

$$p(\mathbf{x}_n | \mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{z}_n, \boldsymbol{\Sigma})$$

$$p(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1 | \boldsymbol{\mu}_0, \mathbf{V}_0)$$

$\mathbf{z}_n$	$=$	$\mathbf{A}\mathbf{z}_{n-1} + \mathbf{w}_n$	$\mathbf{w}$	$\sim \mathcal{N}(\mathbf{w}   \mathbf{0}, \boldsymbol{\Gamma})$
$\mathbf{x}_n$	$=$	$\mathbf{C}\mathbf{z}_n + \mathbf{v}_n$	$\mathbf{v}$	$\sim \mathcal{N}(\mathbf{v}   \mathbf{0}, \boldsymbol{\Sigma})$
$\mathbf{z}_1$	$=$	$\boldsymbol{\mu}_0 + \mathbf{u}$	$\mathbf{u}$	$\sim \mathcal{N}(\mathbf{u}   \mathbf{0}, \mathbf{V}_0)$

# Inference Problem

---

- Finding the marginal distributions for the latent variables conditional on the observation sequence.

$$\widehat{\alpha}(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \mathbf{V}_n)$$

$$\widehat{\beta}(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_{n+1}, \dots, \mathbf{x}_N)$$

$$\gamma(\mathbf{z}_n) = \widehat{\alpha}(\mathbf{z}_n) \widehat{\beta}(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n | \widehat{\boldsymbol{\mu}}_n, \widehat{\mathbf{V}}_n)$$

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = (c_n)^{-1} \widehat{\alpha}(\mathbf{z}_{n-1}) p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n | \mathbf{z}_{-1}) \widehat{\beta}(\mathbf{z}_n)$$

$$c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

---

# Outlines

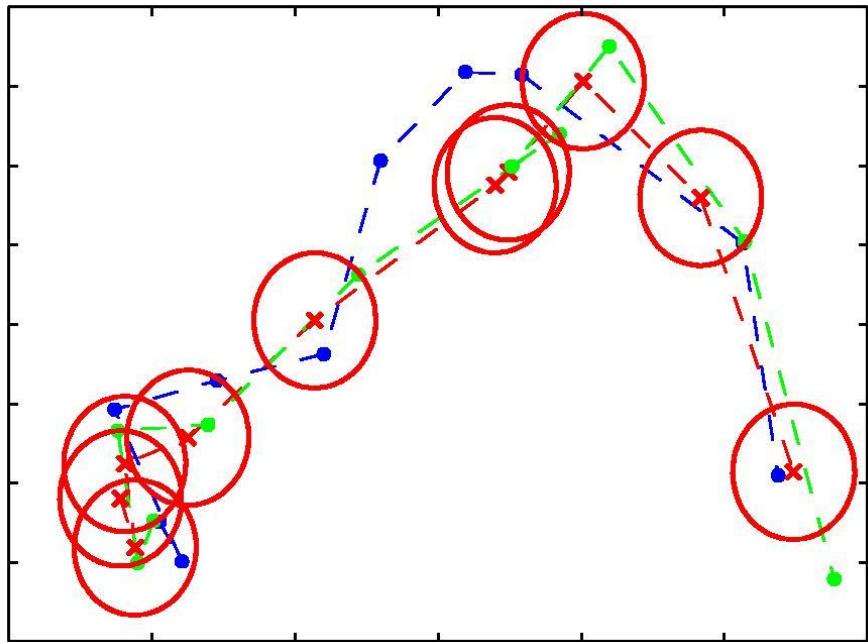
---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-

# Application: Tracking an Moving Object

---

- One of the most important application of the Kalman filter.



An illustration of a linear dynamical system used to track a moving object.

Blue:  $\mathbf{Z}_n$

Green:  $\mathbf{X}_n$

Red:  $\mathbf{Z}_n \mid \mathbf{X}_1, \dots, \mathbf{X}_n$

# Mean and Variance of Kalman Filter

## □ Kalman filter equations

$$\boldsymbol{\mu}_n = \mathbf{A}\boldsymbol{\mu}_{n-1} + \mathbf{K}_n(\mathbf{x}_n - \mathbf{C}\mathbf{A}\boldsymbol{\mu}_{n-1})$$

$$\mathbf{V}_n = (\mathbf{I} - \mathbf{K}_n\mathbf{C})\mathbf{P}_{n-1}$$

$$c_n = \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{A}\boldsymbol{\mu}_{n-1}, \mathbf{C}\mathbf{P}_{n-1}\mathbf{C}^T + \boldsymbol{\Sigma})$$

$$\mathbf{P}_{n-1} = \mathbf{A}\mathbf{V}_{n-1}\mathbf{A}^T + \boldsymbol{\Gamma}$$

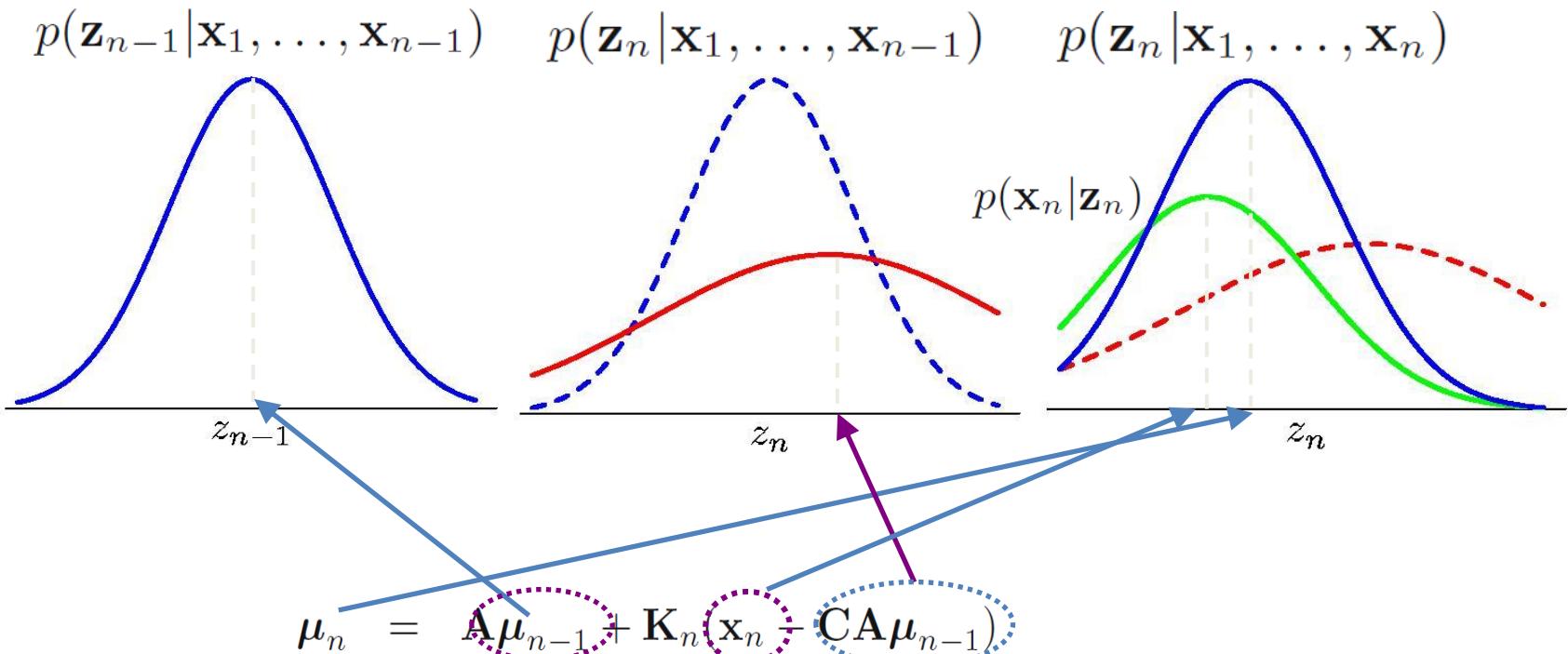
$$\mathbf{K}_n = \mathbf{P}_{n-1}\mathbf{C}^T (\mathbf{C}\mathbf{P}_{n-1}\mathbf{C}^T + \boldsymbol{\Sigma})^{-1}$$

$$\hat{\alpha}(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \mathbf{V}_n)$$

$$c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

# Interpretation of Kalman Filters

- Kalman filter as a process of
  - ✓ Making successive *predictions*
  - ✓ *Correcting* the predictions using new observations.



# Mean and Variance of Full Kalman Filter

---

## □ Full Kalman filter equations

$$\begin{aligned}\hat{\mu}_n &= \mu_n + \mathbf{J}_n (\hat{\mu}_{n+1} - \mathbf{A}\mu_N) \\ \hat{\mathbf{V}}_n &= \mathbf{V}_n + \mathbf{J}_n (\hat{\mathbf{V}}_{n+1} - \mathbf{P}_n) \mathbf{J}_n^T\end{aligned}$$

$$\mathbf{J}_n = \mathbf{V}_n \mathbf{A}^T (\mathbf{P}_n)^{-1}$$

$$\mathbf{A}\mathbf{V}_n = \mathbf{P}_n \mathbf{J}_n^T$$

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_N) = \mathcal{N}(\mathbf{z}_n | \hat{\mu}_n, \hat{\mathbf{V}}_n)$$

# Covariance of Sequential States

---

- Joint posterior of sequential states is Gaussian

$$\begin{aligned}\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) &= p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{z}_{n-1}, \mathbf{z}_n) p(\mathbf{z}_{n-1}, \mathbf{z}_n)}{p(\mathbf{X})} \\ &= \frac{\mathcal{N}(\mathbf{z}_{n-1} | \boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) \mathcal{N}(\mathbf{z}_n | \mathbf{A}\mathbf{z}_{n-1}, \boldsymbol{\Gamma}) \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{z}_n, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{z}_n | \hat{\boldsymbol{\mu}}_n, \hat{\mathbf{V}}_n)}{c_n \hat{\alpha}(\mathbf{z}_n)}\end{aligned}$$

$$\text{cov}[\mathbf{z}_n, \mathbf{z}_{n-1}] = \mathbf{J}_{n-1} \hat{\mathbf{V}}_n$$

# Learning Problem

---

- Determining the parameters  $\vartheta = \{\mathbf{A}, \Gamma, \mathbf{C}, \Sigma, \mu_0, \mathbf{V}_0\}$  using the *EM algorithm*.
- The complete data  $\{\mathbf{X}, \mathbf{Z}\}$  log likelihood function

$$\begin{aligned} \ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) &= \ln p(\mathbf{z}_1|\boldsymbol{\mu}_0, \mathbf{V}_0) + \sum_{n=2}^N \ln p(\mathbf{z}_n|\mathbf{z}_{n-1}, \mathbf{A}, \Gamma) \\ &\quad + \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{z}_n, \mathbf{C}, \Sigma) \end{aligned}$$

---

# Expectation of Latent Variables

---

- The expectation of latent variables

$$\mathbb{E} [\mathbf{z}_n] = \hat{\boldsymbol{\mu}}_n$$

$$\mathbb{E} [\mathbf{z}_n \mathbf{z}_{n-1}^T] = \mathbf{J}_{n-1} \hat{\mathbf{V}}_n + \hat{\boldsymbol{\mu}}_n \hat{\boldsymbol{\mu}}_{n-1}^T$$

$$\mathbb{E} [\mathbf{z}_n \mathbf{z}_n^T] = \hat{\mathbf{V}}_n + \hat{\boldsymbol{\mu}}_n \hat{\boldsymbol{\mu}}_n^T$$

$$\text{cov}[\mathbf{z}_n, \mathbf{z}_{n-1}] = \mathbf{J}_{n-1} \hat{\mathbf{V}}_n$$

# Expectation of Log Likelihood Function

---

- The expectation of the log likelihood function with respect to  $p(\mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}^{\text{old}})$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} [\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})]$$

$$= -\frac{1}{2} \ln |\mathbf{V}_0| - \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} \left[ \frac{1}{2} (\mathbf{z}_1 - \boldsymbol{\mu}_0)^T \mathbf{V}_0^{-1} (\mathbf{z}_1 - \boldsymbol{\mu}_0) \right] + \text{const}$$

$$= -\frac{N-1}{2} \ln |\boldsymbol{\Gamma}| - \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} \left[ \frac{1}{2} \sum_{n=2}^N (\mathbf{z}_n - \mathbf{A}\mathbf{z}_{n-1})^T \boldsymbol{\Gamma}^{-1} (\mathbf{z}_n - \mathbf{A}\mathbf{z}_{n-1}) \right] + \text{const}$$

$$= -\frac{N}{2} \ln |\boldsymbol{\Sigma}| - \mathbb{E}_{\mathbf{Z}|\boldsymbol{\theta}^{\text{old}}} \left[ \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mathbf{C}\mathbf{z}_n)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \mathbf{C}\mathbf{z}_n) \right] + \text{const.}$$

---

# Maximization of LSD Parameters

---

$$\mu_0^{\text{new}} = \mathbb{E}[\mathbf{z}_1]$$

$$\mathbf{V}_0^{\text{new}} = \mathbb{E}[\mathbf{z}_1 \mathbf{z}_1^T] - \mathbb{E}[\mathbf{z}_1] \mathbb{E}[\mathbf{z}_1^T]$$

$$\mathbf{A}^{\text{new}} = \left( \sum_{n=2}^N \mathbb{E} [\mathbf{z}_n \mathbf{z}_{n-1}^T] \right) \left( \sum_{n=2}^N \mathbb{E} [\mathbf{z}_{n-1} \mathbf{z}_{n-1}^T] \right)^{-1}$$

$$\boldsymbol{\Gamma}^{\text{new}} = \frac{1}{N-1} \sum_{n=2}^N \left\{ \mathbb{E} [\mathbf{z}_n \mathbf{z}_n^T] - \mathbf{A}^{\text{new}} \mathbb{E} [\mathbf{z}_{n-1} \mathbf{z}_n^T] - \mathbb{E} [\mathbf{z}_n \mathbf{z}_{n-1}^T] \mathbf{A}^{\text{new}} + \mathbf{A}^{\text{new}} \mathbb{E} [\mathbf{z}_{n-1} \mathbf{z}_{n-1}^T] (\mathbf{A}^{\text{new}})^T \right\}$$

$$\mathbf{C}^{\text{new}} = \left( \sum_{n=1}^N \mathbf{x}_n \mathbb{E} [\mathbf{z}_n^T] \right) \left( \sum_{n=1}^N \mathbb{E} [\mathbf{z}_n \mathbf{z}_n^T] \right)^{-1}$$

$$\boldsymbol{\Sigma}^{\text{new}} = \frac{1}{N} \sum_{n=1}^N \left\{ \mathbf{x}_n \mathbf{x}_n^T - \mathbf{C}^{\text{new}} \mathbb{E} [\mathbf{z}_n] \mathbf{x}_n^T - \mathbf{x}_n \mathbb{E} [\mathbf{z}_n^T] \mathbf{C}^{\text{new}} + \mathbf{C}^{\text{new}} \mathbb{E} [\mathbf{z}_n \mathbf{z}_n^T] \mathbf{C}^{\text{new}} \right\}$$

---

# Extensions of LDS

---

## □ Problem: Beyond the linear-Gaussian assumption.

- ✓ Considerable interest in *extending the basic linear dynamical system* in order to increase its capabilities.
- ✓ Gaussian  $p(\mathbf{z}_n | \mathbf{x}_n)$  – A significant limitation.

## □ Some extensions

- ✓ Gaussian mixture  $p(\mathbf{z}_n)$ .
- ✓ The extended Kalman filter.
- ✓ The switching state space model
- ✓ The switching hidden Markov model

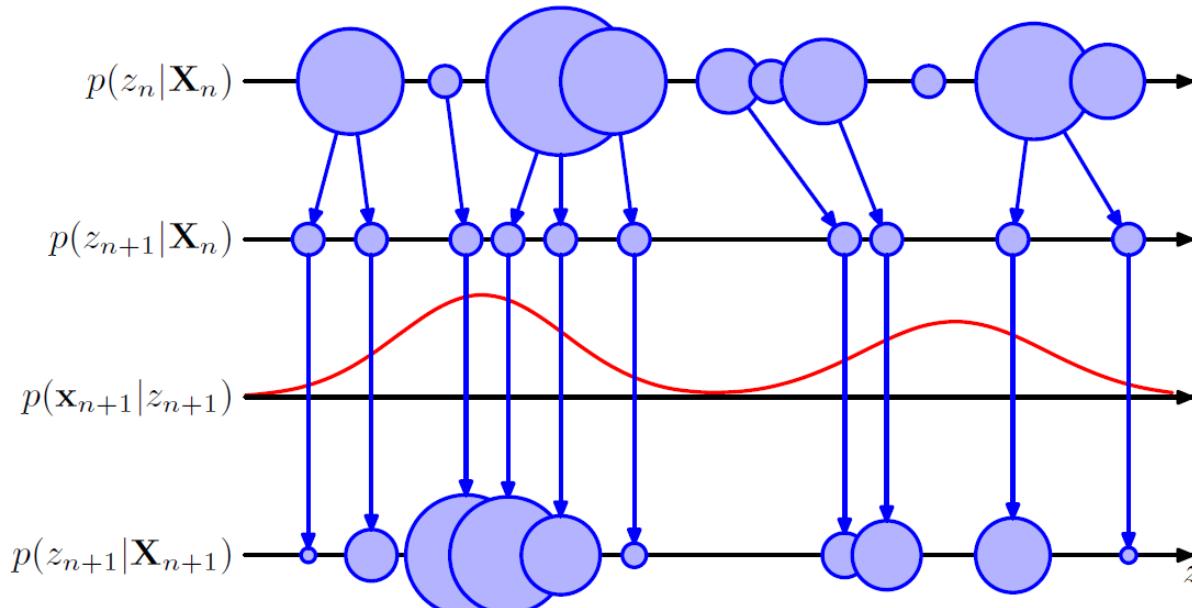
# Particle filters

---

- Non-Gaussian emission density  $p(\mathbf{x}_n | \mathbf{z}_n)$

- ✓ non-Gaussian  $p(z_n | x_1, \dots, x_n)$
- ✓ mathematically intractable integral

- Sampling-importance-resampling

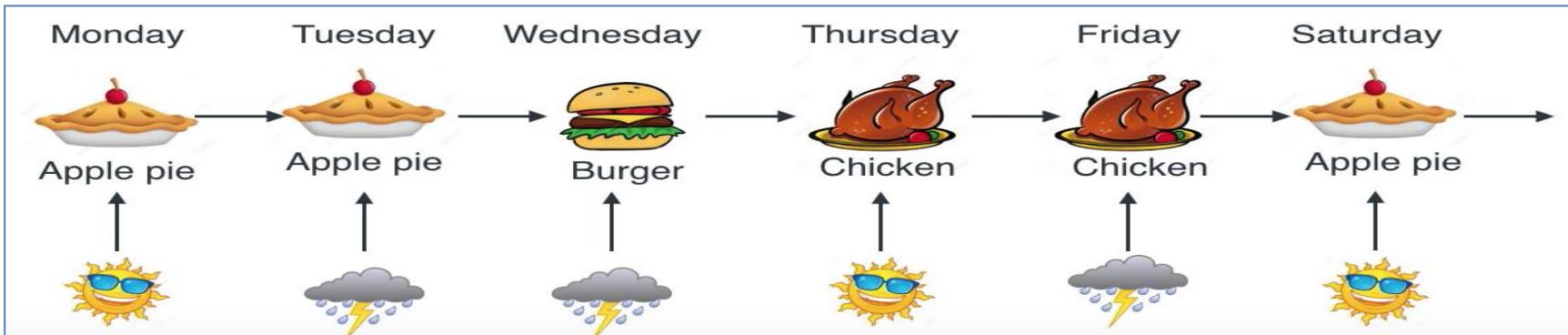


# Outlines

---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-

# Complicated Sequential Data



$$\begin{array}{c}
 \left[ \begin{array}{cccccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{array} \right] \left[ \begin{array}{c} 1 \\ 0 \\ 0 \\ \hline 0 \\ 1 \\ 0 \end{array} \right] = \left[ \begin{array}{c} 1 \\ 0 \\ 0 \\ \hline 0 \\ 1 \\ 0 \end{array} \right] \text{ Apple pie} \\
 \text{Food}
 \end{array}$$

Same      Next day

$$\begin{array}{c}
 \left[ \begin{array}{cccccc} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ \hline 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{array} \right] \left[ \begin{array}{c} 0 \\ 1 \\ 1 \\ \hline 1 \\ 0 \\ 1 \end{array} \right] = \left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ \hline 1 \\ 1 \\ 1 \end{array} \right] \text{ Sun with sunglasses} \\
 \text{Weather}
 \end{array}$$

Same      Next day

$\left[ \begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{array} \right]$

Same

$+$

$\left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{array} \right]$

Same

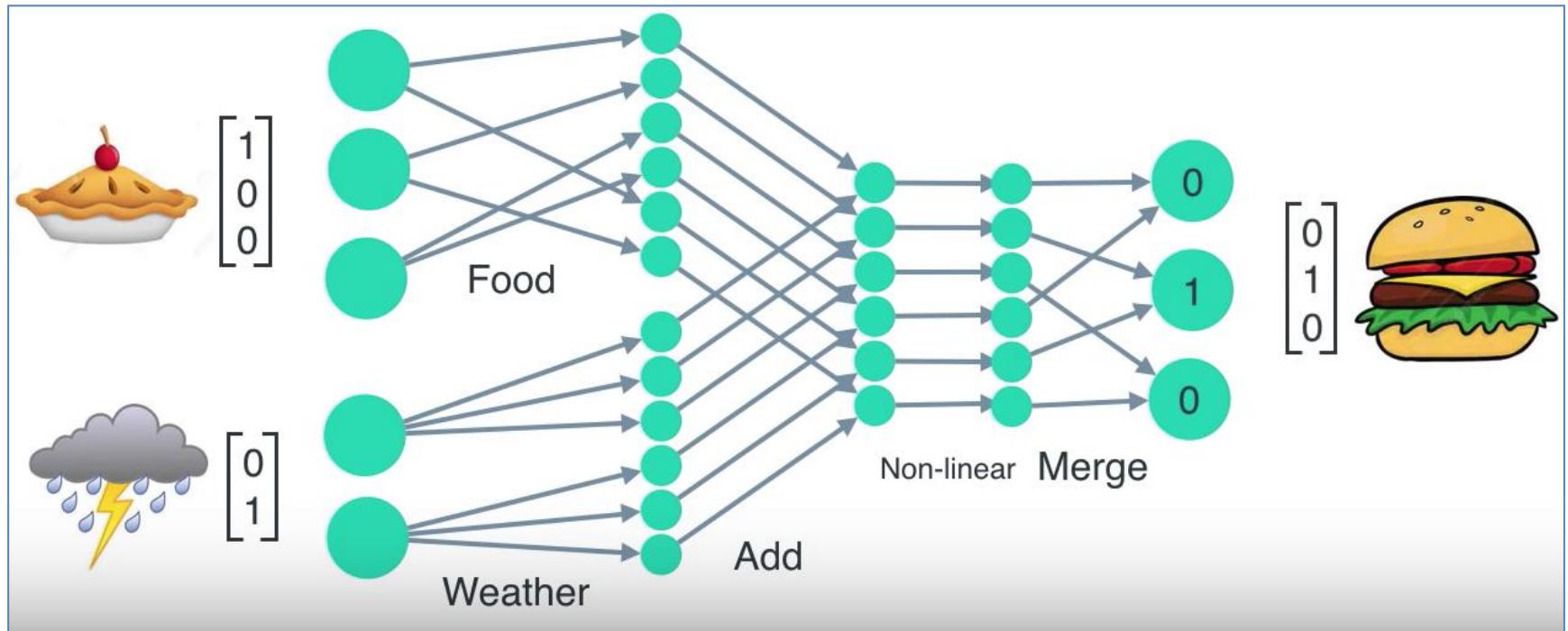
$=$

$\left[ \begin{array}{c} 1 \\ 0 \\ 0 \\ 1 \\ 2 \\ 1 \end{array} \right]$

Next day      Next day

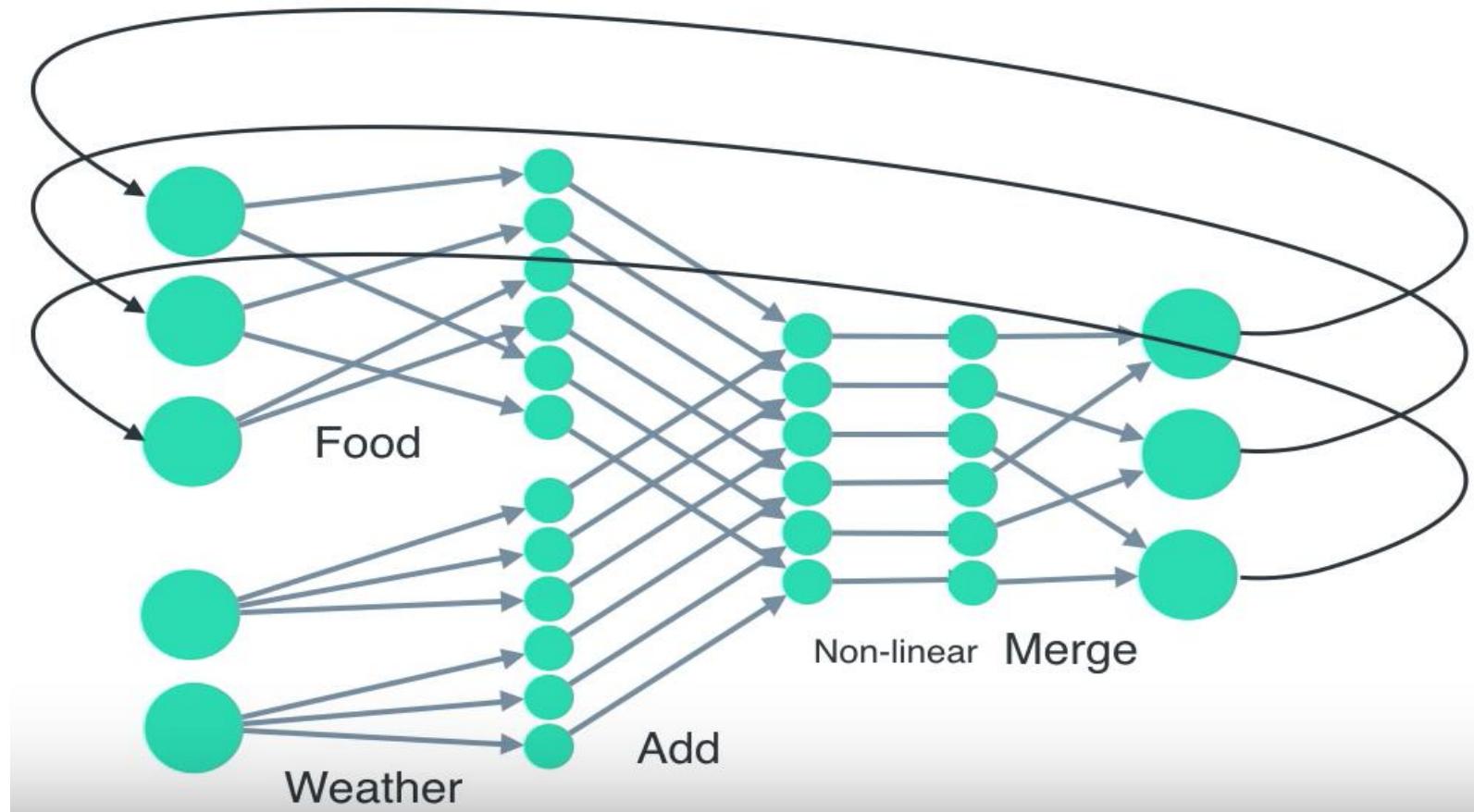
# Neural Network

---



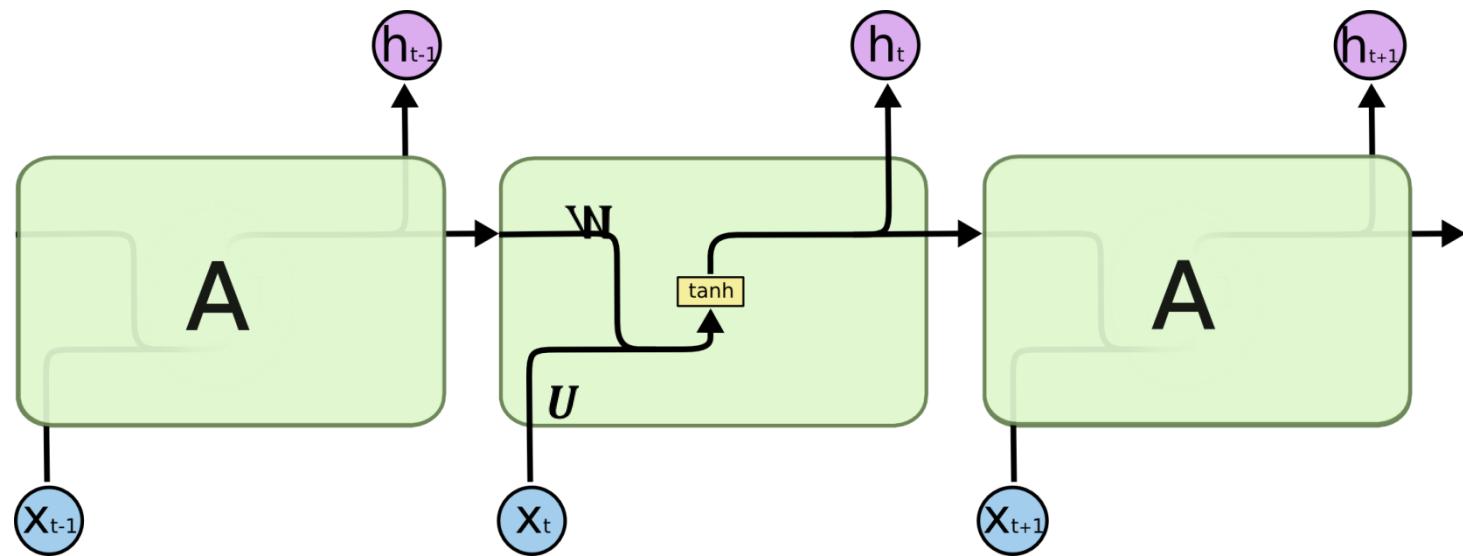
# Recurrent Neural Network

---



# Standard RNN Modules

---



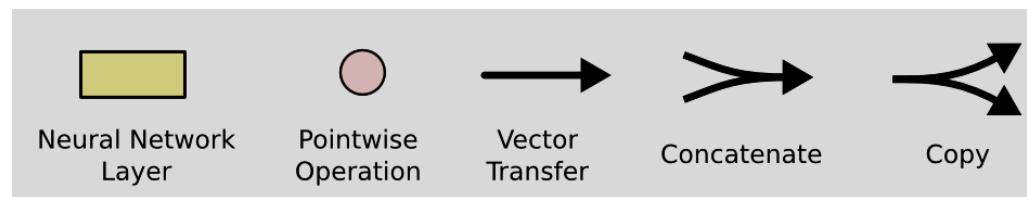
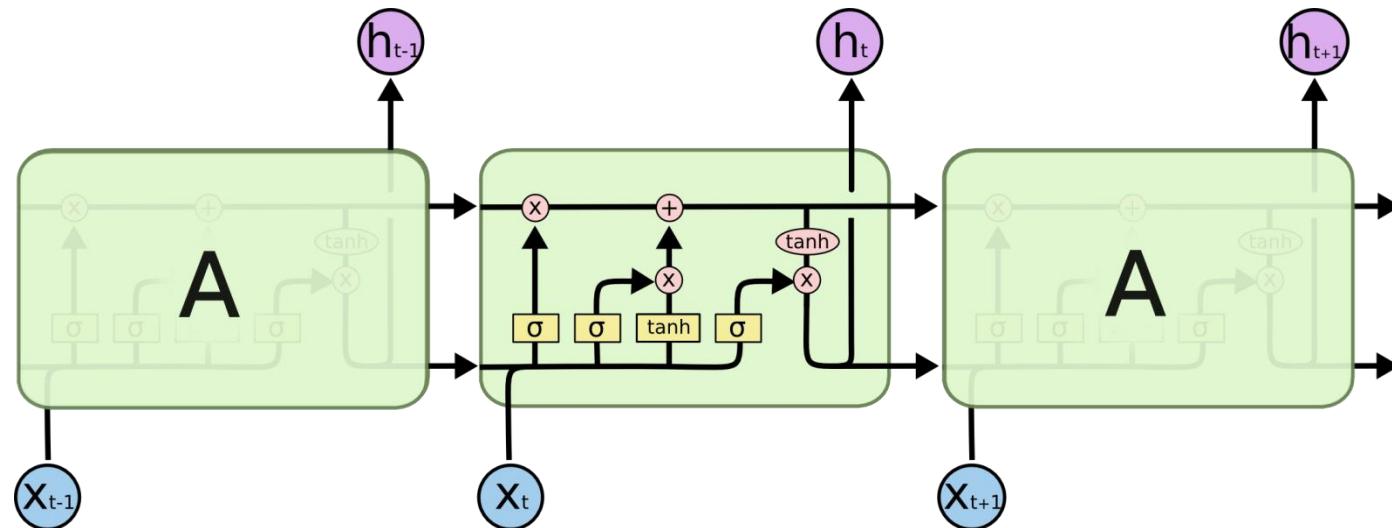
$$\mathbf{h}_t = f(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_{t-1} + \mathbf{b})$$

$\mathbf{h}_t$ : output

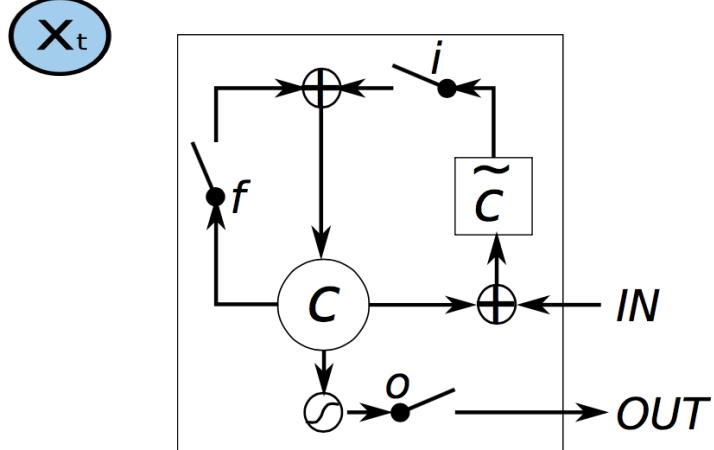
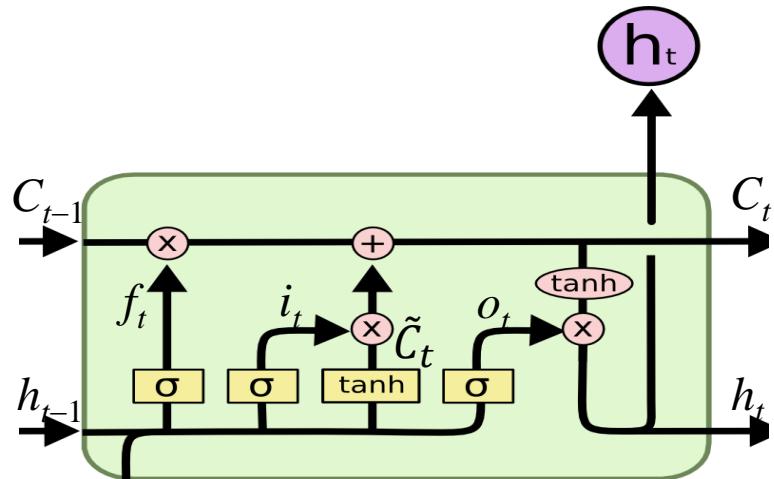
$\mathbf{x}_t$ : input

# Long Short Term Memory

---



# Long Short Term Memory



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

$C_t$ : cell state

$\tilde{C}_t$ : cell state prediction

$f_t$  : forget gate

$i_t$  : input gate

$o_t$  : output gate

$h_t$  : output

$x_t$  : input

# Summary

---

- Hidden Markov Models
  - Maximum Likelihood and EM for HMM
  - Forward-Backward and Sum-Product Algorithms
  - Viterbi and Max-Product Algorithms
  - Linear Dynamics Systems
  - Kalman Filters and LDS Learning
  - RNN and LSTM
-