

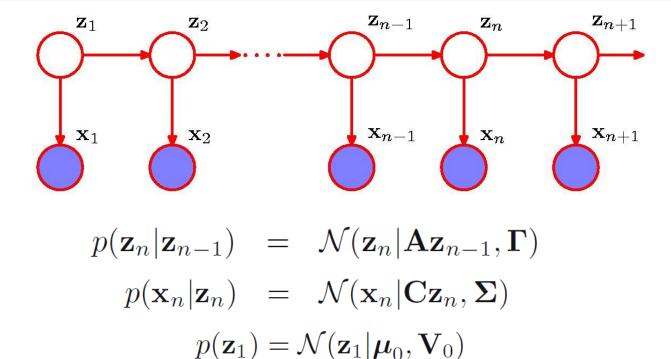
# Learning Objectives

- 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 and Particle Filters
- RNN and LSTM

# Stochastic Linear Dynamical Systems



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

#### Inference Problem

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

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

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

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

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

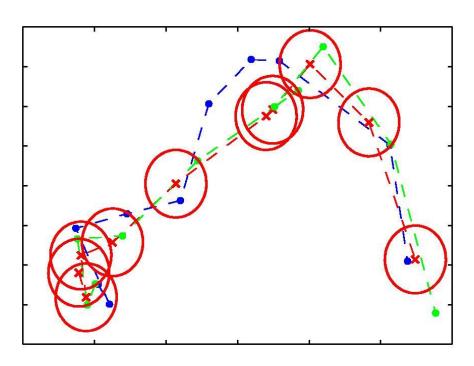
$$c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

#### **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:  $X_n$ 

Red:  $\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_n$ 

#### Mean and Variance of Kalman Filter

■ Kalman filter equations

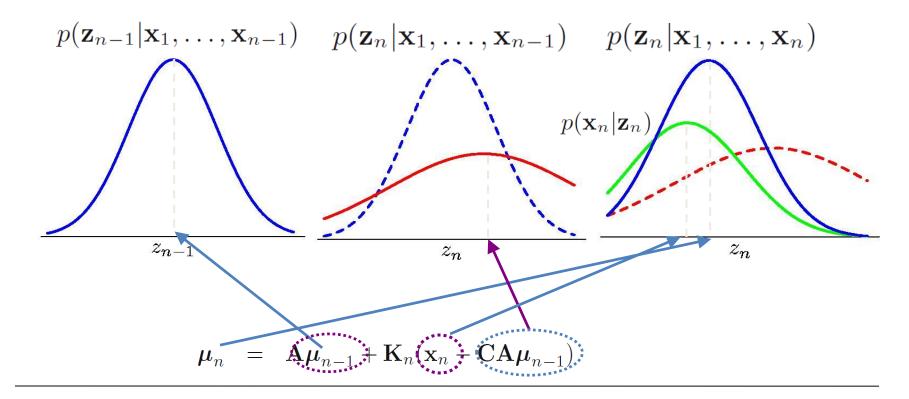
$$egin{array}{lcl} oldsymbol{\mu}_n &=& \mathbf{A}oldsymbol{\mu}_{n-1} + \mathbf{K}_n(\mathbf{x}_n - \mathbf{C}\mathbf{A}oldsymbol{\mu}_{n-1}) \ \mathbf{V}_n &=& (\mathbf{I} - \mathbf{K}_n\mathbf{C})\mathbf{P}_{n-1} \ c_n &=& \mathcal{N}(\mathbf{x}_n|\mathbf{C}\mathbf{A}oldsymbol{\mu}_{n-1},\mathbf{C}\mathbf{P}_{n-1}\mathbf{C}^\mathrm{T} + oldsymbol{\Sigma}) \end{array}$$

$$\mathbf{P}_{n-1} = \mathbf{A} \mathbf{V}_{n-1} \mathbf{A}^{\mathrm{T}} + \mathbf{\Gamma}$$
 $\mathbf{K}_n = \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}} \left( \mathbf{C} \mathbf{P}_{n-1} \mathbf{C}^{\mathrm{T}} + \mathbf{\Sigma} \right)^{-1}$ 

$$\widehat{\alpha}(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \mathcal{N}(\mathbf{z}_n | \boldsymbol{\mu}_n, \mathbf{V}_n)$$
$$c_n = p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1})$$

# 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

$$\widehat{\boldsymbol{\mu}}_n = \boldsymbol{\mu}_n + \mathbf{J}_n \left( \widehat{\boldsymbol{\mu}}_{n+1} - \mathbf{A} \boldsymbol{\mu}_N \right)$$
 $\widehat{\mathbf{V}}_n = \mathbf{V}_n + \mathbf{J}_n \left( \widehat{\mathbf{V}}_{n+1} - \mathbf{P}_n \right) \mathbf{J}_n^{\mathrm{T}}$ 

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

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

# Covariance of Sequential States

☐ Joint posterior of sequential states is Gaussian

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

$$= \frac{\mathcal{N}(\mathbf{z}_{n-1} | \boldsymbol{\mu}_{n-1}, \mathbf{V}_{n-1}) \mathcal{N}(\mathbf{z}_n | \mathbf{A} \mathbf{z}_{n-1}, \boldsymbol{\Gamma}) \mathcal{N}(\mathbf{x}_n | \mathbf{C} \mathbf{z}_n, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{z}_n | \widehat{\boldsymbol{\mu}}_n, \widehat{\mathbf{V}}_n)}{c_n \widehat{\alpha}(\mathbf{z}_n)}$$

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

# **Learning Problem**

- Determining the parameters  $\vartheta = \{A, \Gamma, C, \Sigma, \mu_0, V_0\}$  using the *EM algorithm*.
- $\square$  The complete data  $\{X, Z\}$  log likelihood function

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

## **Expectation of Latent Variables**

■ The expectation of latent variables

$$\mathbb{E} \left[ \mathbf{z}_{n} \right] = \widehat{\boldsymbol{\mu}}_{n}$$

$$\mathbb{E} \left[ \mathbf{z}_{n} \mathbf{z}_{n-1}^{\mathrm{T}} \right] = \mathbf{J}_{n-1} \widehat{\mathbf{V}}_{n} + \widehat{\boldsymbol{\mu}}_{n} \widehat{\boldsymbol{\mu}}_{n-1}^{\mathrm{T}}$$

$$\mathbb{E} \left[ \mathbf{z}_{n} \mathbf{z}_{n}^{\mathrm{T}} \right] = \widehat{\mathbf{V}}_{n} + \widehat{\boldsymbol{\mu}}_{n} \widehat{\boldsymbol{\mu}}_{n}^{\mathrm{T}}$$

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

# Expectation of Log Likelihood Function

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

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

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

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

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

#### Maximization of LSD Parameters

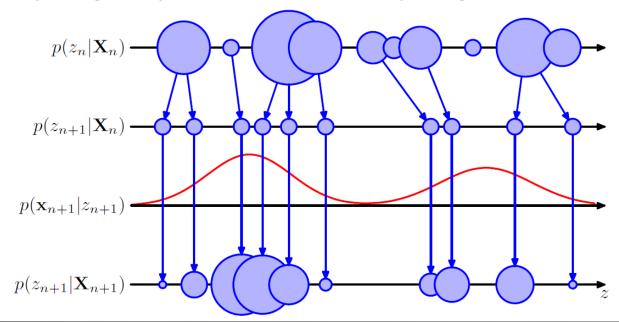
$$\mu_{0}^{\text{new}} = \mathbb{E}[\mathbf{z}_{1}] \\
\mathbf{V}_{0}^{\text{new}} = \mathbb{E}[\mathbf{z}_{1}\mathbf{z}_{1}^{\text{T}}] - \mathbb{E}[\mathbf{z}_{1}]\mathbb{E}[\mathbf{z}_{1}^{\text{T}}] \\
\mathbf{A}^{\text{new}} = \left(\sum_{n=2}^{N} \mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n-1}^{\text{T}}\right]\right) \left(\sum_{n=2}^{N} \mathbb{E}\left[\mathbf{z}_{n-1}\mathbf{z}_{n-1}^{\text{T}}\right]\right)^{-1} \\
\mathbf{\Gamma}^{\text{new}} = \frac{1}{N-1} \sum_{n=2}^{N} \left\{ \mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n}^{\text{T}}\right] - \mathbf{A}^{\text{new}}\mathbb{E}\left[\mathbf{z}_{n-1}\mathbf{z}_{n}^{\text{T}}\right] - \mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n-1}^{\text{T}}\right] \mathbf{A}^{\text{new}} + \mathbf{A}^{\text{new}}\mathbb{E}\left[\mathbf{z}_{n-1}\mathbf{z}_{n-1}^{\text{T}}\right] (\mathbf{A}^{\text{new}})^{\text{T}} \right\} \\
\mathbf{C}^{\text{new}} = \left(\sum_{n=1}^{N} \mathbf{x}_{n}\mathbb{E}\left[\mathbf{z}_{n}^{\text{T}}\right]\right) \left(\sum_{n=1}^{N} \mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n}^{\text{T}}\right]\right)^{-1} \\
\mathbf{\Sigma}^{\text{new}} = \frac{1}{N} \sum_{n=1}^{N} \left\{\mathbf{x}_{n}\mathbf{x}_{n}^{\text{T}} - \mathbf{C}^{\text{new}}\mathbb{E}\left[\mathbf{z}_{n}\right]\mathbf{x}_{n}^{\text{T}} - \mathbf{x}_{n}\mathbb{E}\left[\mathbf{z}_{n}^{\text{T}}\right] \mathbf{C}^{\text{new}} + \mathbf{C}^{\text{new}}\mathbb{E}\left[\mathbf{z}_{n}\mathbf{z}_{n}^{\text{T}}\right] \mathbf{C}^{\text{new}} \right\}$$

#### **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

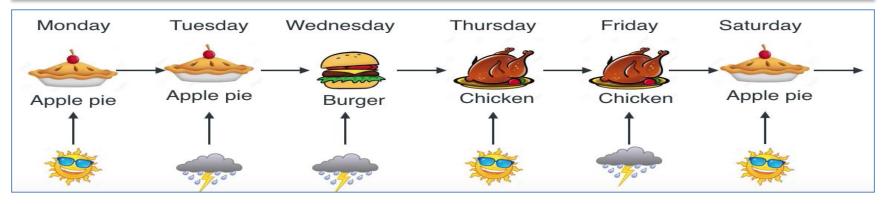
- $\square$  Non-Gaussian emission density  $p(\mathbf{x}_n | \mathbf{z}_n)$ 
  - ✓ non-Gaussian  $p(z_n | x_1, ..., x_n)$
  - ✓ mathematically intractable integral
- Sampling-importance-resampling

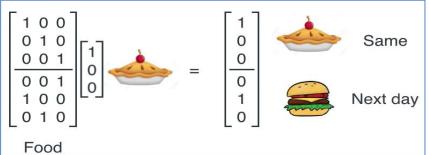


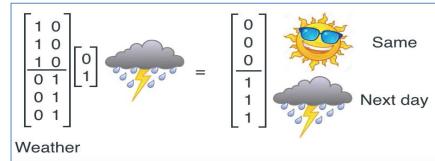
#### **Outlines**

- 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



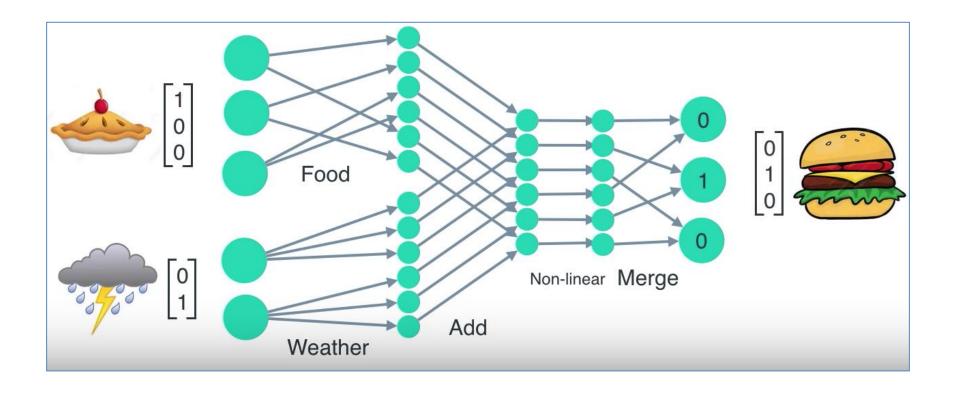




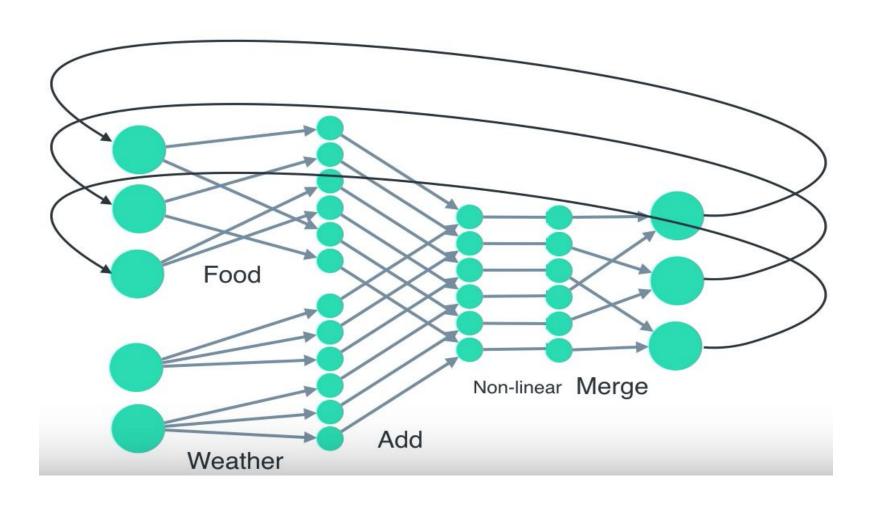


$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$
Same
$$+ \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$
Next day
$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$
Next day

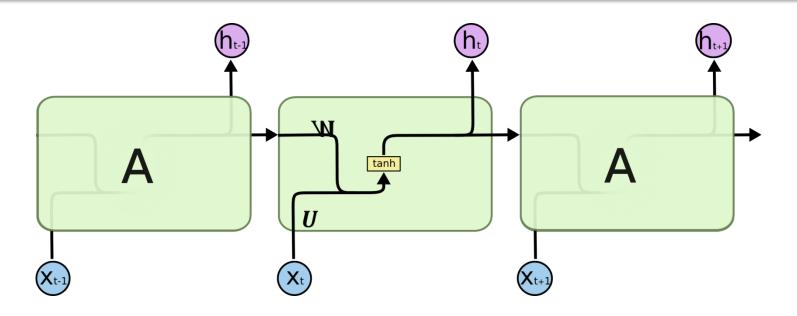
### **Neural Network**



### Recurrent Neural Network



### **Standard RNN Modules**

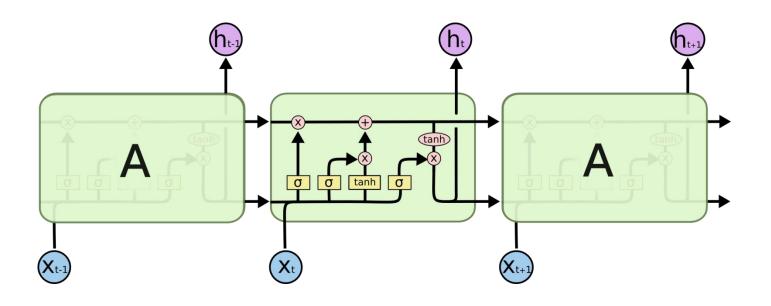


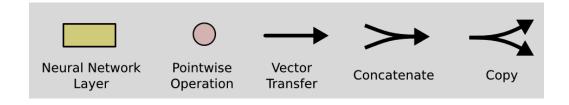
$$h_t = f(Ux_t + Wh_{t-1} + b)$$

 $\mathbf{h}_t$ : output

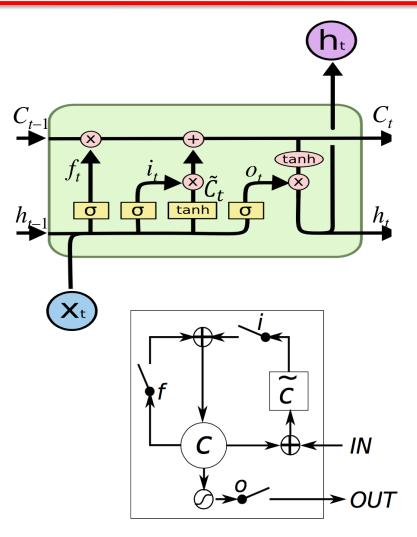
 $x_t$ : input

# Long Short Term Memory





# Long Short Term Memory



$$f_{t} = \sigma (W_{f} \cdot [h_{t-1}, x_{t}] + b_{f})$$

$$i_{t} = \sigma (W_{i} \cdot [h_{t-1}, x_{t}] + b_{i})$$

$$\tilde{C}_{t} = \tanh(W_{C} \cdot [h_{t-1}, x_{t}] + b_{C})$$

$$C_{t} = f_{t} * C_{t-1} + i_{t} * \tilde{C}_{t}$$

$$o_{t} = \sigma (W_{o} [h_{t-1}, x_{t}] + b_{o})$$

$$h_{t} = o_{t} * \tanh(C_{t})$$

 $C_t$ : 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