

D4 and variants

Skorik Sergey

MIPT, 2023

September 23, 2023

1 Backgrounds

2 D4 model

3 D4 model training

4 Discussion

Backgrounds

SSM setup

Consider the problem of modeling time-series data $\mathbf{y}_{1:K}$, $\mathbf{y}_k \in \mathbb{R}^N$, where $k = 1, \dots, K$ using dynamical latent variables $\mathbf{x}_{1:K}$, $\mathbf{x}_k \in \mathbb{R}^M$ with a Markovian property. Under the SSM modeling framework, the joint probability distribution of latent variables and observations can be factorized by conditional probabilities of a generative process defined by

$$\begin{aligned} p(\mathbf{x}_{1:K}, \mathbf{y}_{1:K}) &= p(\mathbf{x}_1, \mathbf{y}_1) \cdot \prod_{t=2}^K p(\mathbf{x}_t, \mathbf{y}_t | \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \\ &= p(\mathbf{x}_1, \mathbf{y}_1) \cdot \prod_{t=2}^K p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \cdot p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) \end{aligned}$$

Backgrounds

SSM setup

Using Markovian property

$$p(\mathbf{x}_t | \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) = p(\mathbf{x}_t | \mathbf{x}_{t-1})$$

$$p(\mathbf{y}_t | \mathbf{x}_t, \mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1}) = p(\mathbf{y}_t | \mathbf{x}_t)$$

Thus, we can write

$$p(\mathbf{x}_{1:K}, \mathbf{y}_{1:K}) = p(\mathbf{x}_1) \cdot p(\mathbf{y}_1 | \mathbf{x}_1) \cdot \prod_{t=2}^K p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{y}_t | \mathbf{x}_t)$$

The posterior distribution is defined by the following recursive solution

$$p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t}) \propto p(\mathbf{y}_t | \mathbf{x}_t) \cdot p(\mathbf{x}_t | \mathbf{x}_{t-1}) \cdot p(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})$$

D4 setup

Setup

Let's assume $\mathbf{h}_t = \mathbf{y}_{1:t-1}$, we can rewrite the posterior distribution of $\mathbf{x}_{1:t}$ given $\mathbf{y}_{1:t}$ as

$$p(\mathbf{x}_{1:t}|\mathbf{y}_t, \mathbf{h}_t) = \frac{p(\mathbf{x}_{1:t}, \mathbf{y}_t, \mathbf{h}_t)}{p(\mathbf{y}_t, \mathbf{h}_t)} = \frac{p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{x}_{1:t-1}, \mathbf{h}_t)p(\mathbf{x}_{1:t-1}, \mathbf{x}_t, \mathbf{h}_t)}{p(\mathbf{y}_t, \mathbf{h}_t)}$$

Using Markovian assumption we can rewrite

$$p(\mathbf{x}_{1:t}|\mathbf{y}_t, \mathbf{h}_t) = \frac{p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{h}_t)p(\mathbf{x}_t, \mathbf{x}_{t-1})p(\mathbf{x}_{1:t-1}, \mathbf{h}_t)p(\mathbf{h}_t)}{p(\mathbf{y}_t, \mathbf{h}_t)}$$

Using Bayesian rule to change $p(\mathbf{y}_t|\mathbf{x}_t, \mathbf{h}_t)$ we can reduce $p(\mathbf{y}_t, \mathbf{h}_t)$ in the numerator and denominator and get

$$p(\mathbf{x}_{1:t}|\mathbf{y}_t, \mathbf{h}_t) = \frac{p(\mathbf{x}_t|\mathbf{y}_t, \mathbf{h}_t)}{p(\mathbf{x}_t|\mathbf{h}_t)}p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{1:t-1}|\mathbf{y}_{t-1}, \mathbf{h}_{t-1}) \quad (1)$$

D4 setup

D4 model

D4 is comprised of two equations

- A state transition equation

$$\mathbf{x}_t | \mathbf{x}_{t-1} \sim g(\mathbf{x}_{t-1}; \boldsymbol{\omega}) \quad (2)$$

- prediction process equation

$$\mathbf{x}_t | \mathbf{y}_t, \mathbf{h}_t \sim f(\mathbf{y}_t, \mathbf{h}_t; \Omega) \quad (3)$$

We need to integrate $p(\mathbf{x}_t | \mathbf{h}_t)$:

$$p(\mathbf{x}_t | \mathbf{h}_t) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{y}_{t-1}, \mathbf{h}_{t-1}) d\mathbf{x}_{t-1}$$

This can be done more efficiently using the sequential (recursive) sampling procedure: **Smoothed sequential importance sampling**

D4 model

Computational complexity

The filtering for SSM requires $\mathcal{O}(NK)$ operations to sample one trajectory of the state approximately distributed according to $p(\mathbf{x}_{1:K}, \mathbf{y}_{1:K})$. On the other hand, the D4 requires $\mathcal{O}(NK^2)$ operations to sample one path for the same distribution. But we can control $\mathbf{h}_t = \mathbf{y}_{1:t-1}$ by replacing with a fixed length history L which can reduce the computational cost to $\mathcal{O}(NKL)$.

D4 model training

EM algorithm

State variables \mathbf{x}_t are not directly observed. For this settings we can use the EM algorithm:

$$\boldsymbol{\theta}^{r+1} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^r) \quad (4)$$

In our case the Q is defined by:

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^r) = \mathbb{E}_{\mathbf{x}_{0:K}|\mathbf{y}_{1:K};\boldsymbol{\theta}^r} [\log (p(\mathbf{Y}, \mathbf{X}|\boldsymbol{\theta}))] \quad (5)$$

where

$$p(\mathbf{Y}, \mathbf{X}|\boldsymbol{\theta}) = p(\omega_0; \mathbf{x}_0) \prod_{t=1}^K p(\mathbf{x}_t|\mathbf{x}_{t-1}; \omega) \prod_{t=1}^K \frac{p(\mathbf{x}_t|\mathbf{y}_t, \mathbf{h}_t; \boldsymbol{\Omega})}{p(\mathbf{x}_t|\mathbf{h}_t; \boldsymbol{\Omega})}$$

and $\boldsymbol{\theta}^r = \{\omega_0^r; \omega^r; \boldsymbol{\Omega}^r\}$. For the simplicity of notation, we use \mathbb{E}_K for $\mathbb{E}_{\mathbf{x}_{0:K}|\mathbf{y}_{1:K};\boldsymbol{\theta}^r}$.

D4 model training

E-step

We need to estimate $Q(\theta|\theta^r)$. We can expand function, as

$$Q(\theta|\theta^r) = \mathbb{E}_K[\log p(\omega_0; \mathbf{x}_0) + \\ + \sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{x}_{t-1}; \omega) + \sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{h}_t; \Omega)] - \mathbb{E}_K \left[\sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{h}_t; \Omega) \right]$$

We can bound the second term

$$\mathbb{E}_K \left[\sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{h}_t; \Omega) \right] = -KL[p(\mathbf{x}_t | \mathbf{y}_{1:K}; \theta^r) || p(\mathbf{x}_t | \mathbf{h}_t; \Omega)] - \\ - \mathbb{H}[p(\mathbf{x}_t | \mathbf{y}_{1:K}; \theta^r)]$$

D4 model training

M-step

The following optimization problem

$$\begin{aligned} \max_{\theta} \mathbb{E}_K [\log p(\omega_0; \mathbf{x}_0) \sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{x}_{t-1}; \omega) + \sum_{t=1}^K \log p(\mathbf{x}_t | \mathbf{y}_t, \mathbf{h}_t; \Omega)] \\ \text{s.t. } \sum_{t=1}^K KL[p(\mathbf{x}_t | \mathbf{y}_{1:K}; \theta^r) || p(\mathbf{x}_t | \mathbf{h}_t; \Omega)] + \mathbb{H}[p(\mathbf{x}_t | \mathbf{y}_{1:K}; \theta^r)] < \varepsilon \end{aligned}$$

D4 training

Algorithm 1 D4 Learning Algorithm

```
1: procedure Regularized-EM-for-D4( $\mathbf{y}_{1:K}, \boldsymbol{\theta}^{(0)}, \lambda, D, L$ )
2:    $\mathbf{h}_k \leftarrow \{\mathbf{y}_{k-L:k-1}\}, Q^0 \leftarrow 0$ 
3:   Do
4:      $Q^{max} \leftarrow Q^r$ 
5:      $\tilde{\mathbf{x}}_{1:K}^{1:D} \leftarrow \mathbf{x}_{1:K}^{1:D}$ 
6:     Sample  $D$  smoothed trajectories using equations 8 and 11,  $\mathbf{x}_{1:K}^{1:D} \sim p_{\boldsymbol{\theta}}(\mathbf{x}_{1:K}^{1:D} \mid \mathbf{y}_{1:K})$ 
7:      $\boldsymbol{\theta}^r, Q^r = \text{Update} - \text{Model}(\mathbf{y}_{1:K}, \tilde{\mathbf{x}}_{1:K}^{1:D}, \mathbf{x}_{1:K}^{1:D}, \boldsymbol{\theta}^{(r-1)}, \lambda)$ 
8:     DoWhile  $\{Q^r > Q^{max}\}$ 
9:     return  $\boldsymbol{\theta}^{(r-1)}, Q^{max}$ 
10: end procedure
11: procedure Update-Model( $\mathbf{y}_{1:K}, \tilde{\mathbf{x}}_{1:K}^{1:D}, \mathbf{x}_{1:K}^{1:D}, \boldsymbol{\theta}^{(r-1)}, \lambda$ )
12:    $\{\boldsymbol{\omega}_0^{(r-1)}, \boldsymbol{\omega}^{(r-1)}, \boldsymbol{\Omega}^{(r-1)}\} = \boldsymbol{\theta}^{(r-1)}$ 
13:   Update  $Q^r$  using equation 17 evaluate at  $\boldsymbol{\theta}^{(r-1)}$ 
14:   Update  $\boldsymbol{\Omega}^r, \boldsymbol{\omega}^r$ , and  $\boldsymbol{\omega}_0^r$  using gradients calculated by equations 18, 26, and 27; respectively
15:   return  $\{\boldsymbol{\omega}_0^r, \boldsymbol{\omega}^r, \boldsymbol{\Omega}^r\}, Q^r$ 
16: end procedure
```

Discussion

- The place of the article in current area
- Match notation in the article and our setup