University of Tehran

Probabilistic Multivariate Time-series Forecasting

Ali Izadi

=

Problem

- Multivariate time-series: $Y \in N \times T$
- Features: $X \in N \times T \times M$

در زمان t فیچر های X را در اختیار داریم و هدف پیش بینی Y است.

Y از یک فضای نهای تولید شده است. بنابراین یک state space در اختیار داریم.

$$egin{aligned} L_1 &\sim N(\mu_1,\, \Sigma_1) \ L_t &= M_t L_{t-1} \,+\, \epsilon_t \ Z_t &= N_t L_t \,+\, \xi_t \ \end{pmatrix} egin{aligned} \epsilon_t &\sim N(0,\, \Sigma_t) \ \xi_t &\sim N(0,\, \Gamma_t) \end{aligned}$$

$$egin{aligned} heta_t \, = \, \mu_1, \, \Sigma_1, \, \{M_t, \, \Sigma_t\}_{t=2}^T, \, \{N_t, \, \Gamma_t\}_{t=1}^T \end{aligned}$$

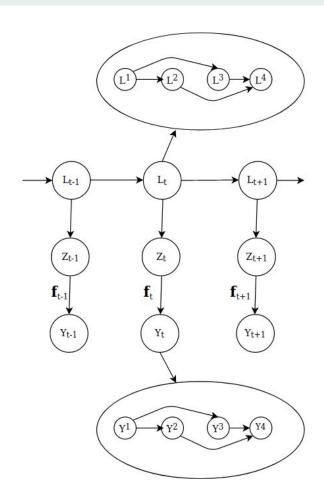
- دیدیم که اگر رابطه بین Z و Y با استفاده از normalizing flow یعنی تابع f غیر خطی شود هم چنان یک state space خطی
 داریم که تنها Yها توسط و ارون تابع f به فرم ساده تری تبدیل میشوند.
 - هم چنین پار امتر های تتا توسط یک RNN که از فیچر ها به دست می آید در زمان t تخمین زده میشوند.

$$egin{aligned} L_1 &\sim N(\mu_1,\,\Sigma_1) \ L_t &= M_t L_{t-1} \,+\, \epsilon_t \ Y_t &= f_t (N_t L_t \,+\, \xi_t) \end{aligned} \qquad egin{aligned} \epsilon_t &\sim N(0,\,\Sigma_t) \ \xi_t &\sim N(0,\,\Gamma_t) \end{aligned}$$

Normalizing Kalman Filter

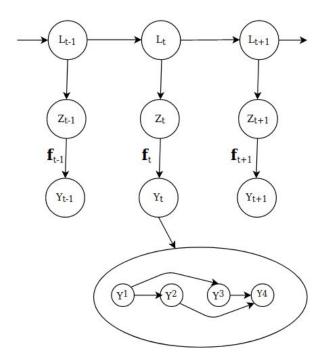
$$egin{aligned} heta_t &= \mu_1, \, \Sigma_1, \, \{M_t, \, \Sigma_t\}_{t=2}^T, \, \{N_t, \, \Gamma_t\}_{t=1}^T \ heta_t &= \sigma(h_t \, = \, \Psi(X_t, \, h_{t-1})) \end{aligned}$$





- دنبال روشی هستیم تا بتواند bayesian network یا dag بین transition بین stateها و هم چنین بین stateها را مشخص کند.
- و همچنان بتوان به از خاصیت kalman filter یعنی فرم بسته filtering و likelihood استفاده کرد.
- روش های موجود یا از قبل bayesian network استفاده را در اختیار دارند یا از روش های variational استفاده اند.





در ابتدا سراغ یادگیری dag بین observationها میرویم.

ایده:

استفاده از ایده مقاله graphical normalizing flow به عنوان یک flow در مقاله normalizing kalman filter تا normalizing kalman filter یک جدید قابلیت یادگیری ساختار در قسمت observation را داشته باشد.



Normalizing Kalman filter

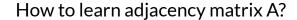
$$egin{aligned} P(Y_{1:T}) \ = \ \prod_{t=1}^T P_{LGM}(Z_t \,|\, Z_{1:t-1}) \,ig| \mathrm{det}\, J_{g(y_t)} ig| \ where \ \ Z_t \ = \ g_t(Y_t) \ \ and \ g_t = f_t^{-1} \end{aligned}$$

and P_{LGM} denotes predictive distribution of kalman filter

 Adding learnable dag A instead of autoregressive normalizing flow

$$P(Y_{1:T}) \ = \ \prod_{t=1}^T P_{LGM}(Z_t \, | \, Z_{1:t-1}) \, \prod_{i=1}^N \left| rac{\mathrm{d} g^iig(Y^i_t, \, h^i(Y_t \circ \, A_{i,:})ig)}{\mathrm{d} Y^i_t}
ight|$$

 $where \circ denotes \, element - wise \, product \, between \, vector \, Y$ $and \, the \, i^{th} \, row \, of \, A$



• Optimization:

$$L(\phi) \, = \, P(Y_{1:T}, \, \phi) \, = \, \prod_{t=1}^T P_{LGM}(Z_t \, | \, Z_{1:t-1}) \, \prod_{i=1}^N \left| rac{\mathrm{d} g^iig(Y^i_t, \, h^i(Y_t \circ \, A_{i,:})ig)}{\mathrm{d} Y^i_t} \,
ight| \, + \, \lambda \, \|A \ \max_{\phi}(L(\phi) \,)$$

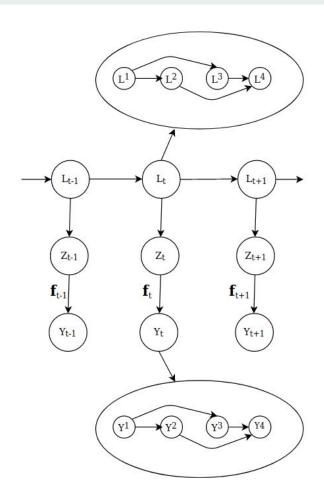
 $egin{aligned} s.\,t. \ h(\phi) \,=\, Tr\,e^{A_\phi}\,-\,N\,=\,0 \end{aligned}$

Augmented Lagrangian

$$egin{split} & \max_{\phi}(L(\phi,\,
ho_t,\,\sigma_t)) \,=\, L(\phi) \,-\,
ho_t h(\phi) \,-\, rac{\sigma_t}{2} h(\phi)^2 \ & \
ho_{t+1} \,\leftarrow\,
ho_t \,+\, \sigma_t h(\phi_t^\cdot) \end{split}$$

$$|\sigma_{t+1} \leftarrow \eta \sigma_t \, if \, h(\phi_t^{\cdot}) \, > \, \gamma hig(\phi_{t-1}^{\cdot}ig) \, else \, \sigma_t$$





 حال به سراغ اضافه کردن dag بین state ها علاوه بر observation ها میرویم.

• ابده

• NO TEARS:
$$F(W) = \ell(W; \mathbf{X}) + \lambda \|W\|_1 = \frac{1}{2n} \|\mathbf{X} - \mathbf{X}W\|_F^2 + \lambda \|W\|_1.$$

• Same adjacency matrix.

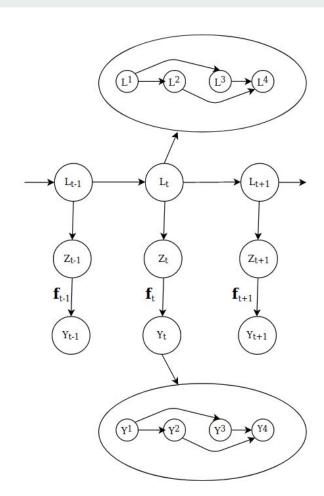
$$egin{aligned} L_t &= M_t L_{t-1} + \epsilon_t & \epsilon_t \sim N(0, \, \Sigma_t) \ \hat{L}_t &= \left(M_t L_{t-1}
ight)^T W \, + \, \epsilon_t \end{aligned}$$

$$egin{align} F(\phi) &= \sum_{t=1}^T \Bigl(L_t \, - \, \hat{L}_t \Bigr) \, + \, \lambda ||W||_1 \ & \ F(\phi) \, = \, \sum_{t=1}^T \Bigl(L_t \, - \, (M_t L_{t-1})^T W \, - \, \epsilon_t \Bigr) \, + \, \lambda ||W||_1 \ & \ \end{array}$$

optimization

$$egin{aligned} F(\phi) &= \sum_{t=1}^T \Bigl(L_t \, - \, (M_t L_{t-1})^T W \, - \, \epsilon_t \Bigr) \, + \, \lambda ||W||_1 \ \min_{\phi} (F(\phi)) \ s. \, t. \ h(\phi) &= Tr \, e^{W_\phi} \, - \, D \, = \, 0 \end{aligned}$$





• پس برای داشتن مدلی مطابق با رو به رو دو optimization خواهیم داشت که میتوان به صورت یک Optimization با دو constraint نوشت که برای این حالت نیز مشابه با زمانی که یک constraint داشتیم روش augmented ا lagrangian با دو constraint استفاده میشود.

https://en.wikipedia.org/wiki/Augmented Lagrangian method

Observations dag

$$egin{aligned} L(\phi) &= P(Y_{1:T},\, \phi) \,=\, \prod_{t=1}^T P_{LGM}(Z_t \,|\, Z_{1:t-1}) \,\prod_{i=1}^N \left| rac{\mathrm{d} g^iig(Y_t^i,\, h^i(Y_t\circ\, A_{i,:})ig)}{\mathrm{d} Y_t^i}
ight| \,+\, \lambda \,\|A\|_1 \ \max_{\phi}(L(\phi)\,) \ s.\,t. \ h(\phi) &= Tr\,e^{A_\phi} \,-\, N \,=\, 0 \end{aligned}$$

States dag

$$egin{aligned} F(\phi) &= \sum_{t=1}^T \Bigl(L_t - (M_t L_{t-1})^T W - \epsilon_t\Bigr) + \lambda ||W||_1 \ \min_{\phi} (F(\phi)) \ s.\, t. \ h(\phi) &= Tr \, e^{W_\phi} - D = 0 \end{aligned}$$

Main optimization

$$egin{aligned} LF &= L(\phi) - F(\phi) \ \max_{\phi}(LF(\phi)) \ s.\,t. \ h_1(\phi) &= Tr\,e^{A_\phi} - N = 0 \ h_2(\phi) &= Tr\,e^{W_\phi} - D = 0 \end{aligned}$$

Augmented Lagrangian

$$egin{split} &\max_{\phi}(LF(\phi,\,
ho_t,\,\sigma_t)) \,=\, LF(\phi) \,-\, \sum_{i=1}^2{(
ho_t)_i h_i(\phi)} \,-\, rac{\sigma_t}{2} \sum_{i=1}^2{h(\phi)^2} \ &(
ho_{t+1})_i \,\leftarrow\, (
ho_t)_i \,+\, \sigma_t h_i(\phi_t^{\cdot}) \ &\sigma_{t+1} \,\leftarrow\, \eta \sigma_t \,if \,h(\phi_t^{\cdot}) \,>\, \gamma hig(\phi_{t-1}^{\cdot}ig) \,else\,\sigma_t \end{split}$$



- 1. Graphical Normalizing flow expressiveness. How can we add multiple flows to get more powerful representation like simple normalizing flow.
- 2. It might help us to introduce new Dag learning approach using the power of normalizing flow.
- 3. Current dag learning methods have O(d^3) computational complexity instead of one approach which is linear some of its computation is approximation.
- 4. O(d^3) is challenging for multivariate time-series with more than 1000 dimensions which is not uncommon

Thank you.

aliizadi2030@gmail.com

