

Latent space Data Assimilation by using Deep Learning

M. Peyron, A. FILLION, S. GÜROL, V. MARCHAIS, S. GRATTON, P. BOUDIER, G. GORET
ANITI/Atos, UFTMIP, CERFACS, ANITI, INPT, NVIDIA, Atos

Email : mathis.peyron@atos.net

Mots Clés : Assimilation de données, Apprentissage profond, espace latent, auto-encodeur, surrogate model, Lorenz 96

Biographie – Diplômé de l'ENSEEIH en 2020 - spécialité *HPC & Big Data* -, j'ai réalisé mon stage de fin d'études au sein de l'ANITI (Artificial and Natural Intelligence Toulouse Institute). Enrichi de cette expérience et dans le *continuum* de mon projet de recherche, je me suis engagé dans une thèse CIFRE avec Atos sur le sujet suivant : "Utilisation conjointe du machine learning et de l'assimilation de données pour un apprentissage plus efficace des dynamiques physiques".

Resumé : Considérons un système dynamique $\mathbf{x} \in \mathbb{R}^n$, un modèle \mathcal{M} (entaché d'une erreur $\boldsymbol{\varepsilon}_k$ à l'instant k) tel que $\mathbf{x}_{k+1} = \mathcal{M}(\mathbf{x}_k) + \boldsymbol{\varepsilon}_k$ et un vecteur d'observations $\mathbf{y} \in \mathbb{R}^p$ lui-même erroné de part une acquisition imparfaite : l'assimilation de données consiste à corriger la prédiction modèle \mathbf{x}_k^f faite à l'instant k par la prise en compte de la donnée observée \mathbf{y}_k . La confiance respective qui peut être faite dans la prédiction \mathbf{x}_k^f et dans l'observation \mathbf{y}_k est centrale pour l'analyse.

La théorie de l'assimilation de données est particulièrement bien adaptée aux phénomènes climatiques pour lesquels d'une part certains processus physiques sont méconnus ou mal représentés, et d'autre part le maillage utilisé lors des simulations ne permet pas de capturer les dynamiques de faible échelle. À ceci, s'ajoute le caractère parfois chaotique de dynamiques propres au phénomène étudié. En météorologie, le système de Lorenz 96 synthétise cette dernière difficulté et fait donc l'objet d'une attention particulière dans la conception d'algorithmes d'assimilation de données.

Bien qu'élégante et au coeur des modélisations actuelles les plus pointues en météorologie, en géosciences ou encore en chimie, l'assimilation de données souffre de deux manquements majeurs : une grande difficulté à traiter des données volumineuses ($n > 10^7, p > 10^6$) et un coût d'application du modèle \mathcal{M} parfois prohibitif. On peut également mentionner l'hypothèse forte selon laquelle les erreurs modèle et d'observation suivent des distributions gaussiennes.

Émergeant depuis une dizaine d'années et aujourd'hui en pleine expansion, l'apprentissage profond répond à ces deux difficultés qui entravent le développement plus efficace de l'assimilation de données. En effet, il rend accessible la recherche d'une représentation latente - sous réserve qu'elle existe - de la dynamique d'intérêt : si les calculs s'avèrent excessivement coûteux dans l'espace physique \mathbb{R}^n , la construction d'un espace latent dans \mathbb{R}^ℓ (avec $\ell \ll n$) est alors légitime voire indispensable. Néanmoins, celle-ci n'est envisageable que si le phénomène physique observé possède une dynamique latente de faible dimension ℓ . Dès lors, le recours aux auto-encodeurs qui bénéficient de transformations non-linéaires offre des performances largement accrues comparativement à celles de l'Analyse en Composantes Principales (ACP).

Une nouvelle assimilation de données plus rapide, dite "latente", est alors réalisable à condition de disposer d'un modèle propagatif latent analogue à \mathcal{M} : ici encore, on peut tirer parti des réseaux de neurones afin d'apprendre une telle propagation latente. Deux bénéfices majeurs sont alors escomptés : réduction du coût calculatoire et de la sollicitation mémoire d'une part ; l'assimilation de données latente est susceptible de produire des résultats plus précis d'autre part.

Cette démarche qui a été la nôtre est une approche nouvelle dans la résolution de problèmes physiques, à la croisée de l'assimilation de données et de l'apprentissage machine. Étant donné

un système dynamique $\mathbf{x} \in \mathbb{R}^n$, nous nous plaçons sous l’hypothèse selon laquelle \mathbf{x} possède une représentation latente $\mathbf{z} \in \mathbb{R}^\ell$ avec $\ell \ll n$. Ainsi la variable \mathbf{z} constitue l’essentiel de l’information contenue dans \mathbf{x} et cette seule connaissance latente suffit à générer une approximation fidèle de la variable physique \mathbf{x} .

Dans ce contexte, il est possible d’entraîner conjointement deux réseaux de neurones, un auto-encodeur et un surrogate : l’objectif du premier est d’engendrer un espace latent de dimension ℓ dans lequel le second sera à même de réaliser la propagation temporelle latente. Si l’on note \mathcal{E} , \mathcal{D} et \mathcal{S} l’encodeur, le décodeur et le surrogate respectivement, alors ce double entraînement est rendu possible au travers d’une fonction de coût adéquate :

$$\mathcal{L}(\mathbf{x}_{k:k+C}) = \frac{1}{C} \sum_{c=1}^C \text{MSE}(\mathcal{D}(\mathcal{E}(\mathbf{x}_{k+c}), \mathbf{x}_{k+c})) + \rho \times \text{MSE}(\mathcal{D}(\mathcal{S}^c(\mathcal{E}(\mathbf{x}_k))), \mathbf{x}_{k+c})$$

où MSE désigne la *Mean Square Error*, C le nombre d’itérations d’application du surrogate \mathcal{S} et ρ représente le poids relatif du second terme dans le coût global. La première MSE contraint la performance de l’auto-encodeur - mesure de la qualité de la reconstruction - alors que la seconde impose l’apprentissage d’une propagation temporelle latente. Celle-ci est itérative afin d’assurer la stabilité du réseau : en effet, il y a un forçage sur la capacité du surrogate à produire les éléments $\mathbf{x}_{k+1}, \mathbf{x}_{k+2}, \dots, \mathbf{x}_{k+C}$ à partir de l’unique donnée \mathbf{x}_k .

Dès lors, on dispose des outils nécessaires à la réalisation d’une assimilation de données latente. L’algorithme retenu est une adaptation de l’ETKF-Q (Ensemble Transform Kalman Filter with Model error) proposé par [1] : cette nouvelle méthode est ici désignée par le sigle ETKF-Q-L. Il s’agit d’un algorithme ensembliste présentant une étape de correction de l’erreur modèle et s’appuyant sur les éléments entraînés \mathcal{E}, \mathcal{D} et \mathcal{S} .

Fortement motivée par les considérations théoriques précédentes, cette nouvelle approche requiert toutefois une justification pratique. Aussi avons-nous considéré un système de Lorenz dit “*augmenté*” construit sur la base des équations de Lorenz en dimension 40 : après avoir généré une telle dynamique de Lorenz sur 500 pas de temps, celle-ci a été multipliée par une matrice orthonormale $O \in \mathbb{R}^{400 \times 40}$. Puis, chaque élément de cette nouvelle matrice s’est vu appliquer une transformation polynomiale de degré 3 inversible. Ce système augmenté possède par construction une représentation latente décrite par les équations de Lorenz 96.

L’approche proposée a été comparée à d’autres méthodes plus simples faisant intervenir une ACP ou une régression linéaire. Les critères de comparaison retenus sont la RMSE - *Root Mean Square Error* - et le temps de calcul. L’approche ETKF-Q-L ressort comme étant la plus performante.

Name	RMSE	Inflation	$\sigma_Q/\sigma_{Q_\epsilon}$	GPU		CPU	
				Avg. Time	Std	Avg. Time	Std
<i>ETKF-Q</i>	0.194	1.12	0.07	16.32s	0.20s	16.23s	0.31s
<i>ETKF-Q-P</i>	0.217	1.08	0.1	13.60s	0.19s	17.12s	0.84s
<i>ETKF-Q-L</i>	0.168	1.004	5.10⁻⁵	6.52s	0.12s	6.89s	0.62s
<i>PCA-S-P</i>	0.383	1.145	0.5	13.22s	0.18s	12.18s	0.31s
<i>PCA-S-L</i>	0.383	1.13	0.6	5.93s	0.10s	5.35s	0.16s
<i>PCA-LinReg-P</i>	0.429	1.24	0.7	12.09s	0.29s	11.62s	0.27s
<i>PCA-LinReg-L</i>	0.426	1.2	0.9	5.02s	0.12s	4.94s	0.15s

Davantage de précisions sur la méthode et les résultats obtenus sont disponibles en ligne [2].

Références

- [1] Anthony Fillion, Marc Bocquet, Serge Gratton, Selime Gürol, and Pavel Sakov. An iterative ensemble kalman smoother in presence of additive model error. *SIAM/ASA Journal on Uncertainty Quantification*, 8(1):198–228, 2020.
- [2] Mathis Peyron, Anthony Fillion, Selime Gürol, Victor Marchais, Serge Gratton, Pierre Boudier, and Gael Goret. Latent space data assimilation by using deep learning, 2021.