

On Multilevel Monte Carlo, Unbiased Gradient Estimation For Deep Latent Variable Models

Simulations et méthodes de Monte Carlo

Nail Khelifa Tom Rossa Axel Pinçon

30 Avril 2024

Ensaie Paris

Introduction

Apports de l'article

Comparaison des estimateurs en termes de Biais et Variance

Application : SGD

Introduction

- Les modèles génératifs sont des modèles de ML qui apprennent à générer des données nouvelles et réalistes à partir de la distribution sous-jacente d'un ensemble de données. Ils sont utilisés dans des domaines comme la génération d'images, de textes, ou de sons.
- Dans notre cas, des variables cachées sont utilisées pour modéliser des structures sous-jacentes non observables, on parle alors de modèles à variables latentes.

- **Observations** : on se donne un modèle paramétrique $\{p_\theta : \theta \in \Theta\}$ et un n -échantillon $\mathbf{x} = \{x^{(i)}\}_{i=1}^n \underset{i.i.d}{\sim} p_{\theta^*}$. On note \mathcal{X} l'espace des observations (ainsi $\mathbf{x} \in \mathcal{X}$).
- **Modèle à variables latentes** : on suppose que les données sont générées par un processus aléatoire impliquant une variable aléatoire continue non observée z :
 1. La valeur $z^{(i)}$ est générée à partir d'une **distribution a priori** $p_{\theta^*}(z)$;
 2. Une valeur $x^{(i)}$ est générée à partir d'une distribution conditionnelle $p_{\theta^*}(\cdot|z)$.
- **Modèles paramétriques** : on suppose que la distribution a priori $p_{\theta^*}(z)$ et la vraisemblance $p_{\theta^*}(\mathbf{x}|z)$ proviennent de familles paramétriques de distributions. **On connaît ainsi l'expression de ces deux distributions.**

En revanche une grande partie de ce processus nous est cachée : les véritables paramètres θ^* ainsi que les valeurs des **variables latentes** $z^{(i)}$ nous sont inconnus.

Motivation

- **Objectif** : Une approche classique pour apprendre θ est de choisir, si celle-ci existe, la valeur de ce paramètre qui maximise la log-vraisemblance marginale de l'échantillon définie par :

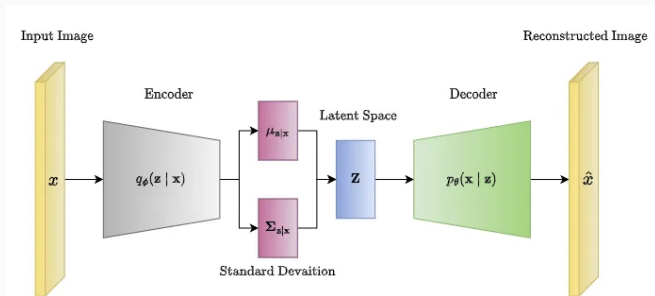
$$\ell(\theta) = \log p_{\theta}(\mathbf{x}) = \log \underbrace{\int_{\mathbf{z}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}}_{\text{Intractable}} \dots$$

- **Problème** : l'intégrale donnée en ci-dessus est intractable
- **Solution** : introduire un **modèle de reconnaissance** (paramétrique) $\{q_{\phi}(\mathbf{z}|\mathbf{x}) : \phi \in \Phi\}$ où $q_{\phi}(\mathbf{z}|\mathbf{x})$ est choisi comme une approximation de la véritable postérieure intractable $p_{\theta}(\mathbf{z}|\mathbf{x})$.

$$\dots = \log \int_{\mathbf{z}} \frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} q_{\phi}(\mathbf{z}|\mathbf{x}) d\mathbf{z} = \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Un exemple de Modèle Génératif: le VAE

Un VAE utilise deux réseaux neuronaux : un encodeur pour encoder des données dans un espace latent, et un décodeur pour générer des données à partir de cet espace.



Apports de l'article

On observe la décomposition suivante :

$$\log p_{\theta}(\mathbf{x}) = \mathcal{L}(\theta, \phi) + \underbrace{\text{KL}(q_{\phi}(z|\mathbf{x}) \parallel p_{\theta}(z|\mathbf{x}))}_{\geq 0} \geq \mathcal{L}(\theta, \phi)$$

où :

- \mathcal{L} est l'ELBO (Evidence Lower Bound).
- La KL Divergence, $\text{KL}(q_{\phi}(z|\mathbf{x}) \parallel p_{\theta}(z|\mathbf{x}))$, mesure la dissimilarité entre la distribution approximative $q_{\phi}(z|\mathbf{x})$ et la distribution vraie $p_{\theta}(z|\mathbf{x})$.

Une première approche : l'IWAE (1)

- **Définition - Importance Weighted Autoencoder Estimator :**

$$\ell_{\text{IAWE}}^{(k)}(\theta, \phi) = \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_k \sim_{i.i.d} q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{1}{k} \sum_{i=1}^k \frac{p_\theta(\mathbf{x}, \mathbf{z}_i)}{q_\phi(\mathbf{z}_i|\mathbf{x})} \right]$$

où $\mathbf{z}_1, \dots, \mathbf{z}_k$ sont échantillonnés indépendamment du modèle de reconnaissance.

- **Inconvénient** : à moins que $q_\phi(\mathbf{z}|\mathbf{x})$ corresponde exactement à $p_\theta(\mathbf{z}|\mathbf{x})$, alors en général $\nabla_{\theta} \ell_{\text{IAWE}}^{(k)}(\theta, \phi) \neq \nabla_{\theta} \ell(\theta)$

$\Rightarrow \ell_{\text{IAWE}}^{(k)}(\theta, \phi)$ est un estimateur consistant mais **biaisé** pour tout k fini.

Une première approche : l'IWAE (2)

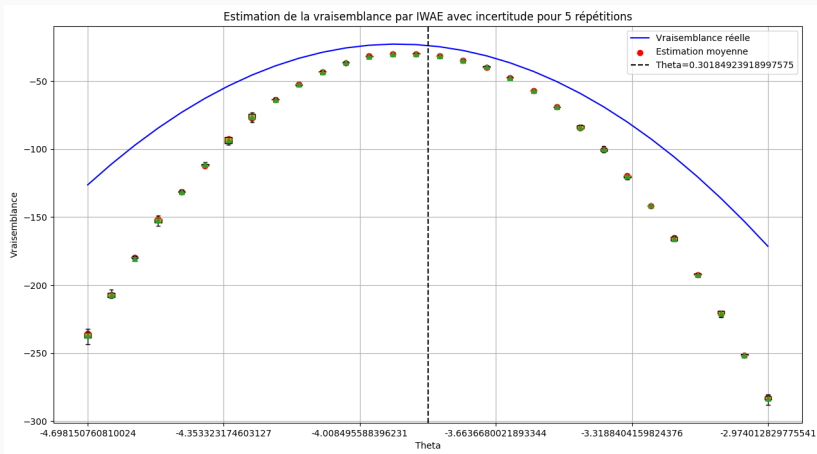


Figure 1: Estimation de la log-vraisemblance par IWAE avec incertitude pour 5 répétitions

Une première approche : l'IWAE (3)

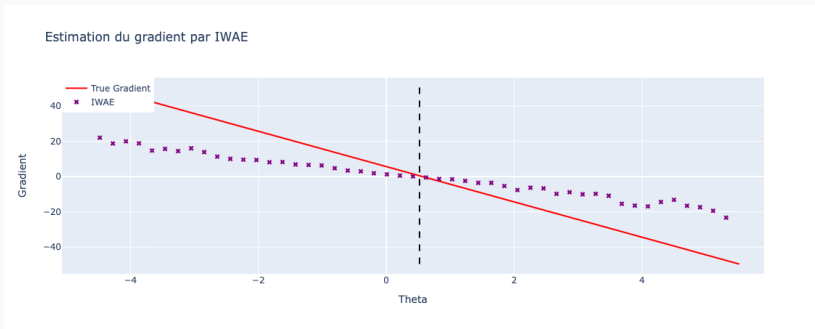


Figure 2: Estimation du gradient par IWAE

Débiaiser un estimateur

- **Idée** : proposer des estimateurs non-biaisés de la vraisemblance et du gradient
- **Proposition** : Désignons une quantité d'intérêt par $l_\infty = \log p_\theta(\mathbf{x})$. Supposons que l_∞ puisse être écrite comme

$$l_\infty = \mathbb{E}[l_0] + \sum_{k=0}^{\infty} \mathbb{E}[\Delta_k]$$

pour les variables aléatoires l_0 et $(\Delta_k)_{k \geq 0}$. Nous pouvons estimer l_∞ de manière non biaisée via les estimateurs suivants ss ou rr :

$$ss = l_0 + \frac{\Delta_K}{p(K)}, \quad rr = l_0 + \sum_{k=0}^K \frac{\Delta_k}{\mathbb{P}(K \geq k)}$$

où $K \sim \text{Geom}(r)$.

Un estimateur non-biaisé : SUMO (1)

- **Définition :**

$$\begin{aligned}\Delta_k^{\text{SUMO}} &:= \hat{\ell}^{(k+2)}(\boldsymbol{\theta}) - \hat{\ell}^{(k+1)}(\boldsymbol{\theta}) \\ &:= \log \left(\frac{1}{k+2} \sum_{i=1}^{k+2} \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right) - \log \left(\frac{1}{k+1} \sum_{i=1}^{k+1} \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right),\end{aligned}$$

et alors :

$$\hat{\ell}^{\text{SUMO}}(\boldsymbol{\theta}) := l_0 + \sum_{k=0}^K \frac{\Delta_k^{\text{SUMO}}}{P(\mathcal{K} \geq k)}$$

- **Inconvénients :** variance potentiellement non bornée

Un estimateur non-biaisé : SUMO (2)

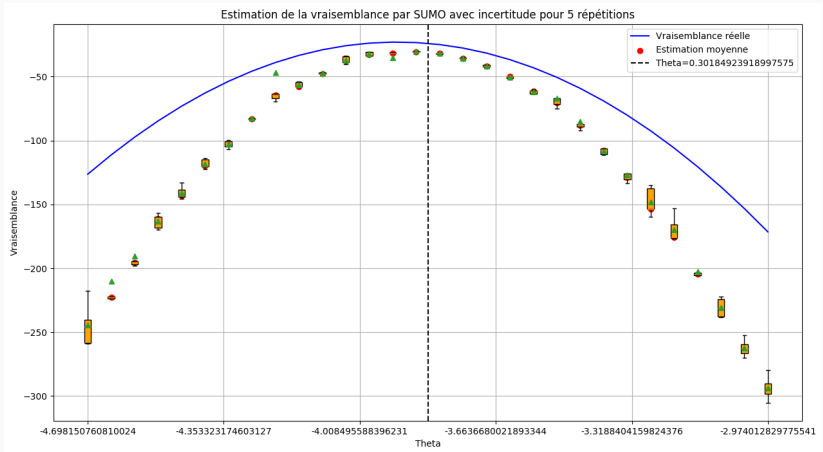


Figure 3: Estimation de la vraisemblance par SUMO

Un estimateur non-biaisé : SUMO (3)

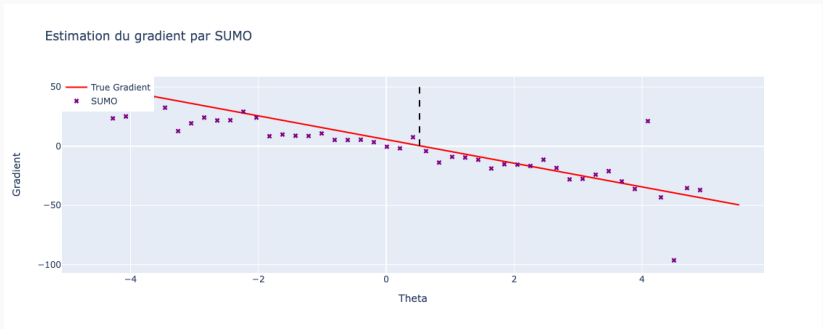


Figure 4: Estimation du gradient par SUMO

- **Motivation** : construire un estimateur non biaisé et computationnellement calculable du gradient pour l'**optimisation SGD**.
- **Problème** : Estimateurs du gradient soit **biaisés** soit **de variance non bornée**.
- **Objectif** : Combiner les estimateurs de la roulette russe aux méthodes Multi Level Monte Carlo pour construire deux estimateurs (ML-SS et ML-RR) non biaisés et de variance bornée du gradient.

Les estimateurs ML-SS et ML-RR (1)

- **Idée** : même idée que pour SUMO, repose sur un choix astucieux de Δ_k
- **Définition** : on pose

$$\Delta_k^{\text{ML}} = \hat{\ell}_{O \cup E}^{(2^{k+1})}(\theta) - \frac{1}{2} \left(\hat{\ell}_O^{(2^k)}(\theta) + \hat{\ell}_E^{(2^k)}(\theta) \right),$$

où \mathbf{z}_i^O , \mathbf{z}_i^E sont deux séquences indépendantes d'échantillons i.i.d. de q_ϕ , où O , E désignent respectivement impair (odd) et pair (even). On a alors

$$\hat{\ell}^{\text{ML-SS}}(\theta) = l_0 + \frac{\Delta_{\mathcal{K}}^{\text{ML}}}{p(\mathcal{K})}, \quad \hat{\ell}^{\text{ML-RR}}(\theta) = l_0 + \sum_{k=0}^K \frac{\Delta_k^{\text{ML}}}{\mathbb{P}(\mathcal{K} \geq k)},$$

Les estimateurs ML-SS et ML-RR (2)

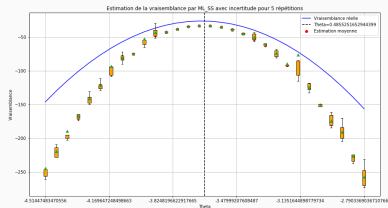


Figure 5: Estimation de la vraisemblance par ML-SS

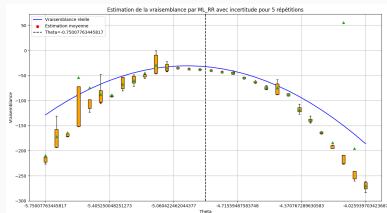


Figure 6: Estimation de la vraisemblance par ML-SS

Les estimateurs ML-SS et ML-RR (3)

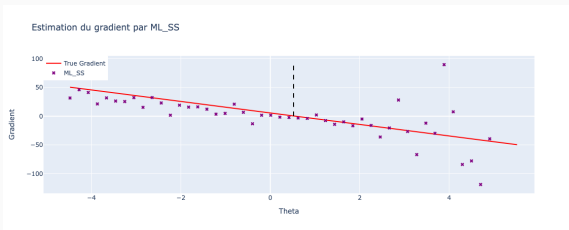


Figure 7: Estimation du gradient par ML-SS

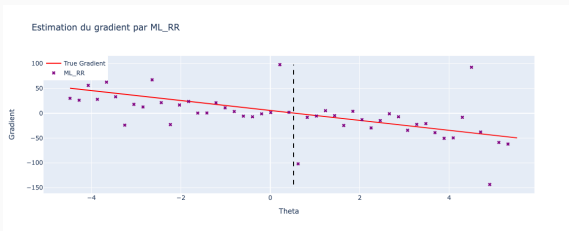


Figure 8: Estimation du gradient par ML-RR

Comparaison des estimateurs en termes de Biais et Variance

Analyse Biais / Variance (1)

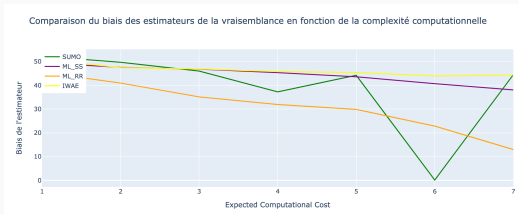


Figure 9: Biais au carré des estimateurs de la log-vraisemblance

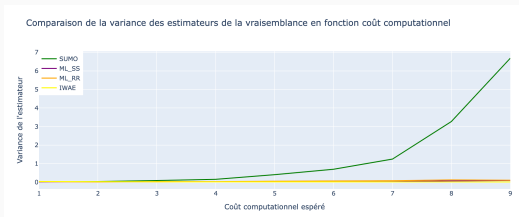


Figure 10: Variance des estimateurs de la log-vraisemblance

Analyse Biais / Variance (2)

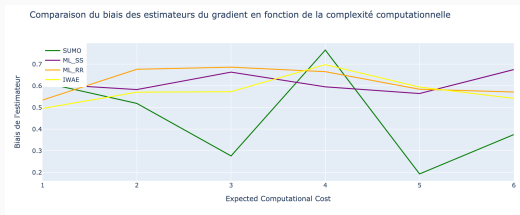


Figure 11: Biais au carré des estimateurs du gradient

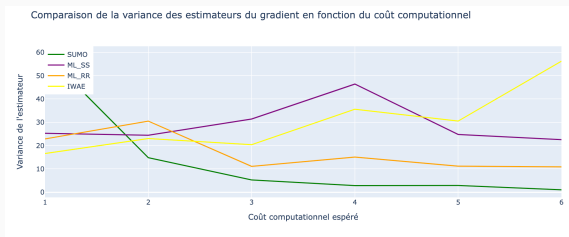


Figure 12: Variance des estimateurs du gradient

Application : SGD

Cadre général de la Descente de Gradient (SGD)

La descente de gradient est un algorithme d'optimisation permet minimiser une fonction convexe. On l'initialise avec un $\theta^{(0)}$ puis :

À chaque itération t , les paramètres θ sont mis à jour en suivant la direction opposée au gradient de notre fonction qui se trouve être la direction de plus forte pente.

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_{\theta} \log p_{\theta^{(t)}}(x), \quad (1)$$

On répète (1) jusqu'à ce qu'un critère d'arrêt prédéfini :

$$\|\theta^{(t+1)} - \theta^{(t)}\| < \epsilon$$

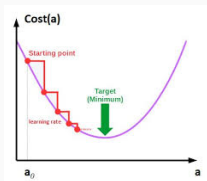


Figure 13: Convergence d'une descente de gradient

Stochastic Gradient Descent (SGD) dans les VAE

A partir de nos différents estimateurs, nous avons effectué une SGD afin de maximiser la log-vraisemblance des données observées :

$$\ell(\theta) = \log p_{\theta}(\mathbf{x}) = \sum_{i=1}^{20} \log p_{\theta}(x^{(i)})$$

Implémentation de la SGD avec la méthode SUMO: à chaque itération t on actualise $\theta^{(t)}$ en ne calculant qu'une composante du gradient

$$\theta^{(t+1)} = \theta^{(t)} + \eta \nabla_{\theta} \hat{\ell}_{(i)}^{\text{SUMO}}(\theta), \quad (2)$$

Où $\hat{\ell}_{(i)}^{\text{SUMO}}(\theta)$ est un estimateur de $\nabla_{\theta} \log p_{\theta}(x^{(i)})$ et $x^{(i)}$ tirée aléatoirement.

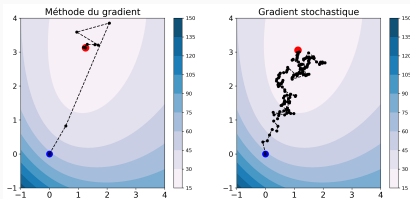


Figure 14: GD vs SGD

Comparaison des estimateurs

Comparaison de la convergence de la SGD avec les différents estimateurs du gradient

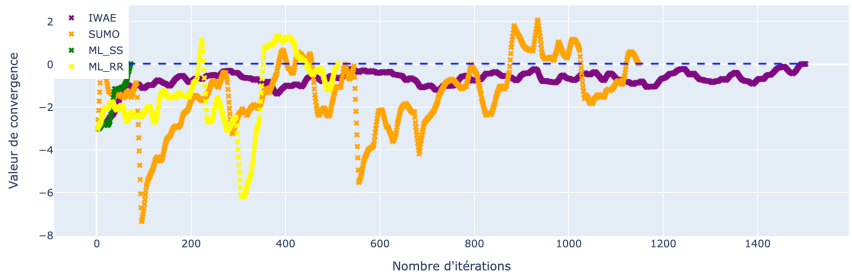


Figure 15: Convergence d'une descente de gradient

Questions?

**Nous vous invitons à demander quelque
éclaircissement que ce soit.**