

VARIATIONAL AUTOENCODERS (VAE)

Jérémie Cabessa

Laboratoire DAVID, UVSQ

INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un espace de dimension plus petite – l'**espace latent (latent space)** – et être capable de les reconstruire ces data.
- ▶ On avait vu que l'encodeur pouvait être utilisé comme un **data compressor**...
- ▶ et le **décodeur** comme un **data generator**.
- ▶ Pour cela, il suffisait de sampler des points dans l'espace latent et de les décoder.

INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un espace de dimension plus petite – l'**espace latent (latent space)** – et être capable de les reconstruire ces data.
- ▶ On avait vu que l'**encodeur** pouvait être utilisé comme un **data compressor...**
- ▶ et le **décodeur** comme un **data generator**.
- ▶ Pour cela, il suffisait de sampler des points dans l'espace latent et de les décoder.

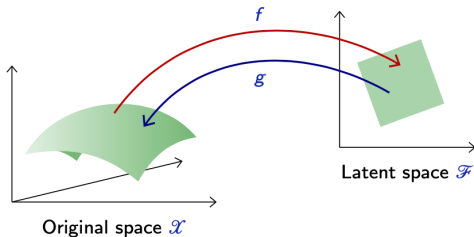
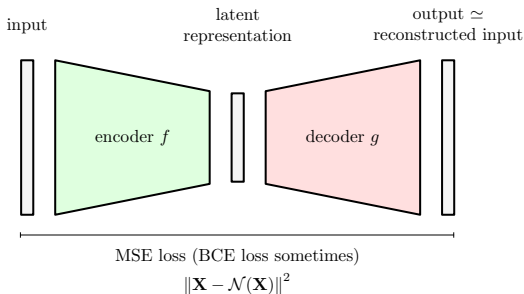
INTRODUCTION

- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un espace de dimension plus petite – l'**espace latent (latent space)** – et être capable de les reconstruire ces data.
- ▶ On avait vu que l'**encodeur** pouvait être utilisé comme un **data compressor...**
- ▶ et le **décodeur** comme un **data generator**.
- ▶ Pour cela, il suffisait de sampler des points dans l'espace latent et de les décoder.

INTRODUCTION

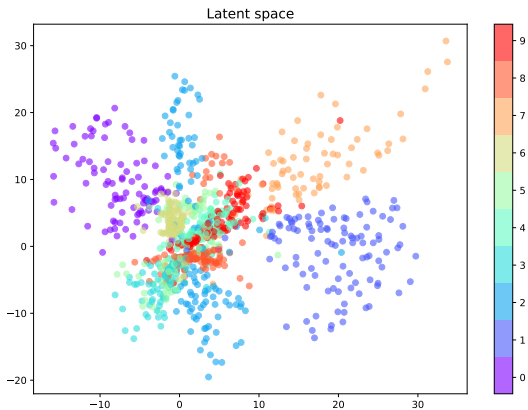
- ▶ **But d'un autoencodeur (rappel):** projeter les data dans un espace de dimension plus petite – l'**espace latent (latent space)** – et être capable de les reconstruire ces data.
- ▶ On avait vu que l'**encodeur** pouvait être utilisé comme un **data compressor**...
- ▶ et le **décodeur** comme un **data generator**.
- ▶ Pour cela, il suffisait de sampler des points dans l'espace latent et de les décoder.

INTRODUCTION



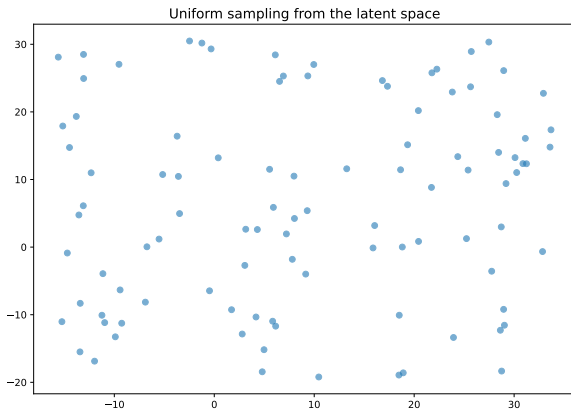
MOTIVATION

- Dans notre cas, on avait un espace latent de dimension 2, et lorsqu'on samplait dans cet espace, les data générées étaient d'assez bonne qualité (à l'oeil).



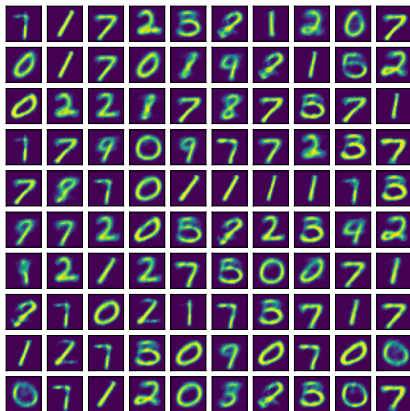
MOTIVATION

- Dans notre cas, on avait un espace latent de dimension 2, et lorsqu'on samplait dans cet espace, les data générées étaient d'assez bonne qualité (à l'oeil).



MOTIVATION

- Dans notre cas, on avait un espace latent de dimension 2, et lorsqu'on samplait dans cet espace, les data générées étaient d'assez bonne qualité (à l'oeil).



MOTIVATION

- ▶ Mais cette situation n'est pas représentative...
- ▶ En pratique, on considère des espace latent de dimensions plus grandes que 2 (on arrive rarement à compresser un signal en dimension 2).
- ▶ Dans ce cas, le sampling dans l'espace latent génère des data beaucoup plus dégradées...

MOTIVATION

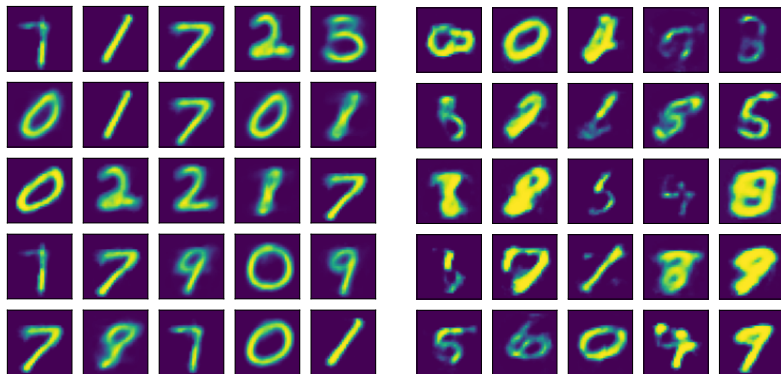
- ▶ Mais cette situation n'est pas représentative...
- ▶ En pratique, on considère des espace latent de dimensions plus grandes que 2 (on arrive rarement à compresser un signal en dimension 2).
- ▶ Dans ce cas, le sampling dans l'espace latent génère des data beaucoup plus dégradées...

MOTIVATION

- ▶ Mais cette situation n'est pas représentative...
- ▶ En pratique, on considère des espace latent de dimensions plus grandes que 2 (on arrive rarement à compresser un signal en dimension 2).
- ▶ Dans ce cas, le sampling dans l'espace latent génère des data beaucoup plus dégradées...

MOTIVATION

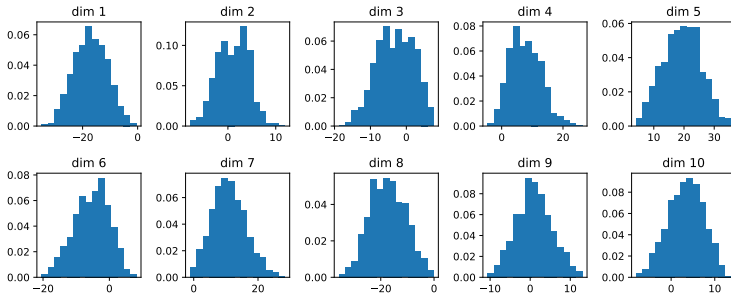
- ▶ Par exemple, si on répète l'opération avec un espace latent de dimension 10, et qu'on sample dans cet espace, les data générées sont bien moins bonnes.



Data generated from latent spaces of dim 2 and 10, respectively.

MOTIVATION

- De plus, la distribution des data dans l'espace latent est très déséquilibrée (non-centered, skewed, etc.).



Distribution of data in the latent space for each dimension.

MOTIVATION

- ▶ Pour un autoencodeur classique, qu'est-ce qui rend le processus de sampling de l'espace latent difficile?
 - ★ Distribution de l'espace latent non-continue: ne se voit pas forcément en 2D, mais en dimension supérieure, il y aura sûrement des "trous".
 - ★ Distribution de l'espace latent non-équilibrée et non-centrée.
- ▶ L'encodeur variationnel (variational autoencoder VAE) vise à pallier ces problèmes.
- ▶ Il force l'espace latent à suivre une certaine distribution, normale en général.

MOTIVATION

- ▶ Pour un autoencodeur classique, qu'est-ce qui rend le processus de sampling de l'espace latent difficile?
- ★ Distribution de l'espace latent non-continue: ne se voit pas forcément en 2D, mais en dimension supérieure, il y aura sûrement des "trous".
- ★ Distribution de l'espace latent non-équilibrée et non-centrée.
- ▶ L'encodeur variationnel (variational autoencoder VAE) vise à pallier ces problèmes.
- ▶ Il force l'espace latent à suivre une certaine distribution, normale en général.

MOTIVATION

- ▶ Pour un autoencodeur classique, qu'est-ce qui rend le processus de sampling de l'espace latent difficile?
- ★ Distribution de l'espace latent non-continue: ne se voit pas forcément en 2D, mais en dimension supérieure, il y aura sûrement des "trous".
- ★ Distribution de l'espace latent non-équilibrée et non-centrée.
- ▶ L'encodeur variationnel (variational autoencoder VAE) vise à pallier ces problèmes.
- ▶ Il force l'espace latent à suivre une certaine distribution, normale en général.

MOTIVATION

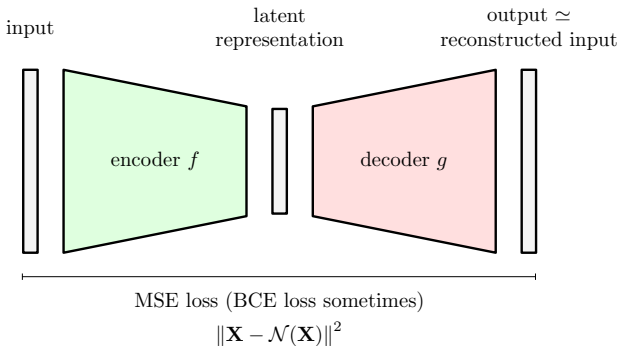
- ▶ Pour un autoencodeur classique, qu'est-ce qui rend le processus de sampling de l'espace latent difficile?
- ★ Distribution de l'espace latent non-continue: ne se voit pas forcément en 2D, mais en dimension supérieure, il y aura sûrement des "trous".
- ★ Distribution de l'espace latent non-équilibrée et non-centrée.
- ▶ L'encodeur **variationnel** (variational autoencoder **VAE**) vise à pallier ces problèmes.
- ▶ Il force l'espace latent à suivre une certaine distribution, normale en général.

MOTIVATION

- ▶ Pour un autoencodeur classique, qu'est-ce qui rend le processus de sampling de l'espace latent difficile?
- ★ Distribution de l'espace latent non-continue: ne se voit pas forcément en 2D, mais en dimension supérieure, il y aura sûrement des "trous".
- ★ Distribution de l'espace latent non-équilibrée et non-centrée.
- ▶ L'encodeur **variationnel** (**variational autoencoder VAE**) vise à pallier ces problèmes.
- ▶ Il force l'espace latent à suivre une certaine distribution, normale en général.

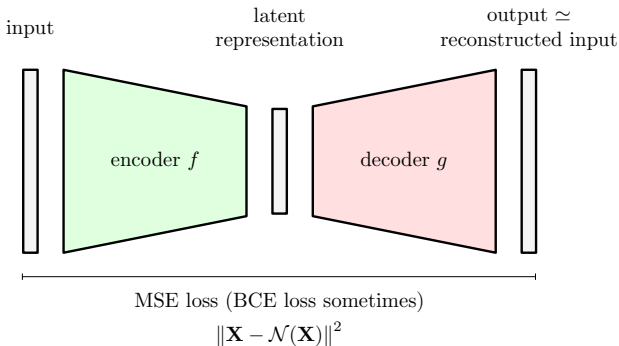
VARIATIONAL AUTOENCODER (VAE)

- **Idée générale:** Imposer une distribution sur l'espace latent et entraîner un décodeur de manière à reconstruire les data.
- On passe d'un bottleneck *déterministe* à un bottleneck *stochastique*.



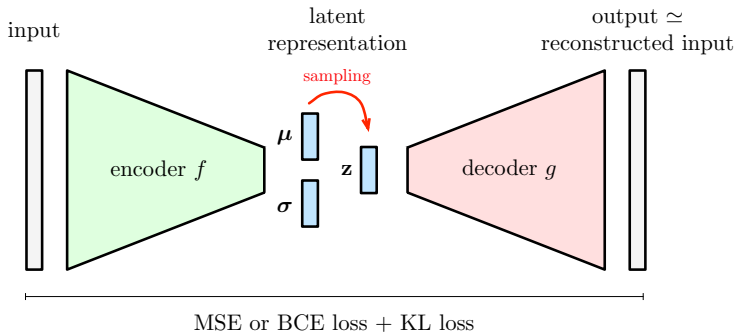
VARIATIONAL AUTOENCODER (VAE)

- **Idée générale:** Imposer une distribution sur l'espace latent et entraîner un décodeur de manière à reconstruire les data.
- On passe d'un bottleneck *déterministe* à un bottleneck *stochastique*.



VARIATIONAL AUTOENCODER (VAE)

- **Idée générale:** Imposer une distribution sur l'espace latent et entraîner un décodeur de manière à reconstruire les data.
- On passe d'un bottleneck *déterministe* à un bottleneck *stochastique*.



VARIATIONAL AUTOENCODER (VAE)

- ▶ Régularisation des data: on impose des contraintes de *continuité* et de *complétude* sur l'espace latent.
- ▶ **Contrainte de continuité:** on aimerait que des points proches de l'espace latent génère des data similaires.
- On introduit un processus de décodage stochastique lors de l'entraînement: chaque data est échantillonnée dans l'espace latent (voisinage) avant d'être décodée.
- ▶ **Contrainte de complétude:** on aimerait que tous les points de l'espace latent génèrent des data qui soient valables.
- Pour "ramasser les data" (éviter les trous), on force les data de l'espace latent à suivre une loi normale centrée réduite $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ Régularisation des data: on impose des contraintes de *continuité* et de *complétude* sur l'espace latent.
 - ▶ **Contrainte de continuité:** on aimerait que des points proches de l'espace latent génère des data similaires.
- On introduit un processus de décodage stochastique lors de l'entraînement: chaque data est échantillonnée dans l'espace latent (voisinage) avant d'être décodée.
- ▶ **Contrainte de complétude:** on aimerait que tous les points de l'espace latent génèrent des data qui soient valables.
- Pour "ramasser les data" (éviter les trous), on force les data de l'espace latent à suivre une loi normale centrée réduite $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ Régularisation des data: on impose des contraintes de *continuité* et de *complétude* sur l'espace latent.
- ▶ **Contrainte de continuité:** on aimerait que des points proches de l'espace latent génère des data similaires.
- On introduit un processus de décodage stochastique lors de l'entraînement: chaque data est échantillonnée dans l'espace latent (voisinage) avant d'être décodée.
- ▶ **Contrainte de complétude:** on aimerait que tous les points de l'espace latent génèrent des data qui soient valables.
- Pour "ramasser les data" (éviter les trous), on force les data de l'espace latent à suivre une loi normale centrée réduite $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

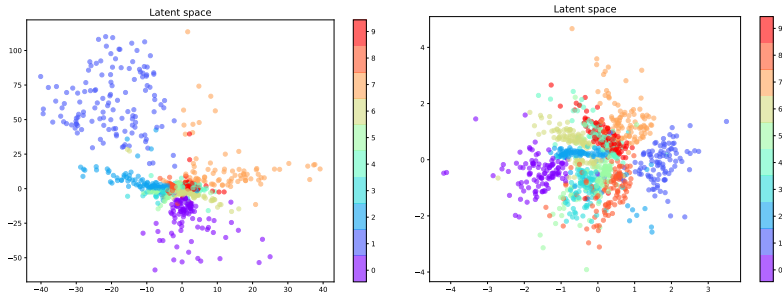
VARIATIONAL AUTOENCODER (VAE)

- ▶ Régularisation des data: on impose des contraintes de *continuité* et de *complétude* sur l'espace latent.
- ▶ **Contrainte de continuité:** on aimerait que des points proches de l'espace latent génère des data similaires.
- On introduit un processus de décodage stochastique lors de l'entraînement: chaque data est échantillonnée dans l'espace latent (voisinage) avant d'être décodée.
- ▶ **Contrainte de complétude:** on aimerait que tous les points de l'espace latent génèrent des data qui soient valables.
- Pour “ramasser les data” (éviter les trous), on force les data de l'espace latent à suivre une loi normale centrée réduite $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

VARIATIONAL AUTOENCODER (VAE)

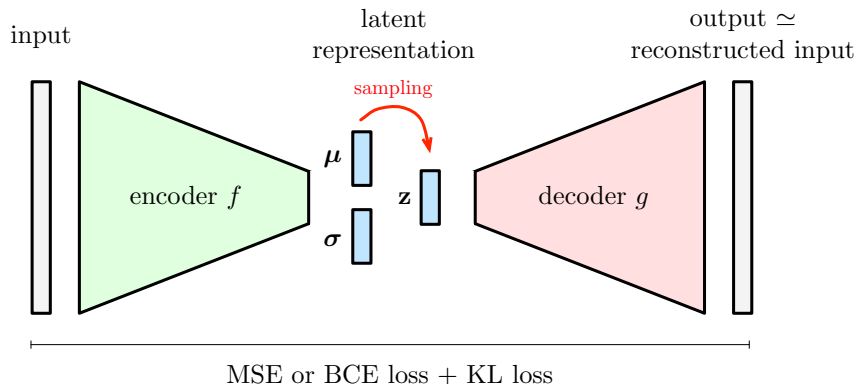
- ▶ Régularisation des data: on impose des contraintes de *continuité* et de *complétude* sur l'espace latent.
- ▶ **Contrainte de continuité:** on aimerait que des points proches de l'espace latent génère des data similaires.
- On introduit un processus de décodage stochastique lors de l'entraînement: chaque data est échantillonnée dans l'espace latent (voisinage) avant d'être décodée.
- ▶ **Contrainte de complétude:** on aimerait que tous les points de l'espace latent génèrent des data qui soient valables.
- Pour “ramasser les data” (éviter les trous), on force les data de l'espace latent à suivre une loi normale centrée réduite $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

VARIATIONAL AUTOENCODER (VAE)

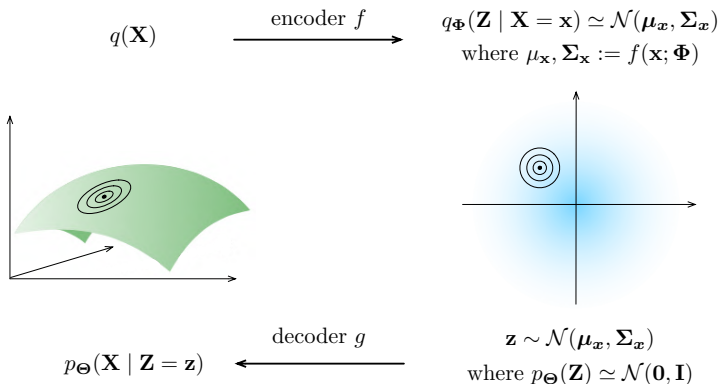


Left: less continuous and complete latent space; **Right:** more continuous and complete latent space.

VARIATIONAL AUTOENCODER (VAE)

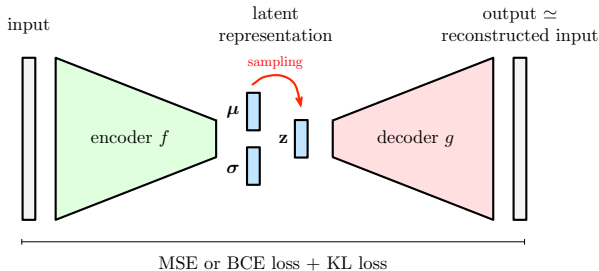


VARIATIONAL AUTOENCODER (VAE)



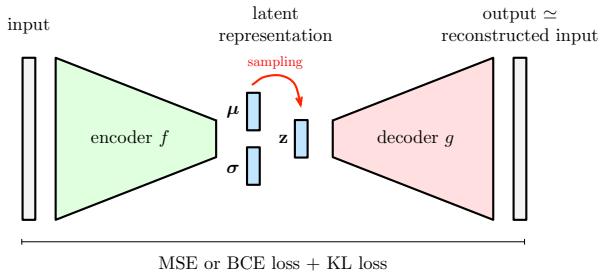
VARIATIONAL AUTOENCODER (VAE)

- ▶ L'encodeur est un deep neural network $f(\cdot; \Phi)$.
- ▶ Le sampler $s(\cdot; \mu, \Sigma)$ sample selon une loi normale $\mathcal{N}(\mu, \Sigma)$.
- ▶ Le décodeur est un deep neural network $g(\cdot; \Theta)$.
- ▶ La composition de $f(\cdot; \Phi)$, $s(\cdot; \mu, \Sigma)$ et $g(\cdot; \Theta)$ forme un réseau de neurones stochastique $\mathcal{N}(\cdot; \Phi, \Theta)$.



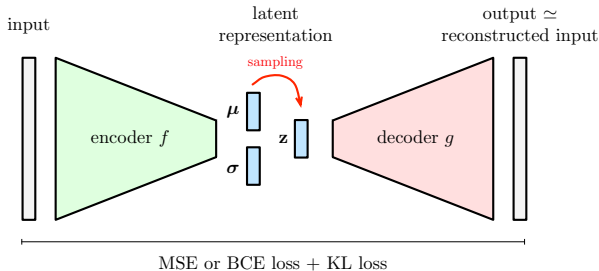
VARIATIONAL AUTOENCODER (VAE)

- ▶ L'**encodeur** est un deep neural network $f(\cdot; \Phi)$.
- ▶ Le **sampler** $s(\cdot; \mu, \Sigma)$ sample selon une loi normale $\mathcal{N}(\mu, \Sigma)$.
- ▶ Le **décodeur** est un deep neural network $g(\cdot; \Theta)$.
- ▶ La composition de $f(\cdot; \Phi)$, $s(\cdot; \mu, \Sigma)$ et $g(\cdot; \Theta)$ forme un **réseau de neurones stochastique** $\mathcal{N}(\cdot; \Phi, \Theta)$.



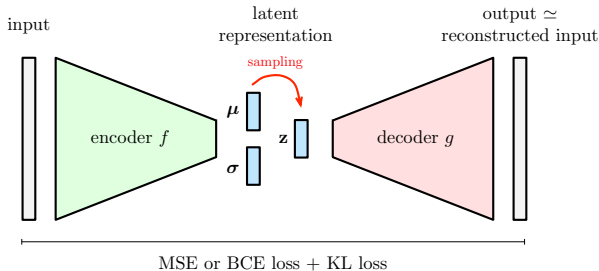
VARIATIONAL AUTOENCODER (VAE)

- ▶ L'**encodeur** est un deep neural network $f(\cdot; \Phi)$.
- ▶ Le **sampler** $s(\cdot; \mu, \Sigma)$ sample selon une loi normale $\mathcal{N}(\mu, \Sigma)$.
- ▶ Le **décodeur** est un deep neural network $g(\cdot; \Theta)$.
- ▶ La composition de $f(\cdot; \Phi)$, $s(\cdot; \mu, \Sigma)$ et $g(\cdot; \Theta)$ forme un **réseau de neurones stochastique** $\mathcal{N}(\cdot; \Phi, \Theta)$.



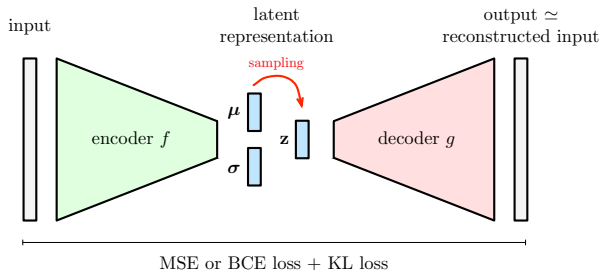
VARIATIONAL AUTOENCODER (VAE)

- ▶ L'**encodeur** est un deep neural network $f(\cdot; \Phi)$.
- ▶ Le **sampler** $s(\cdot; \mu, \Sigma)$ sample selon une loi normale $\mathcal{N}(\mu, \Sigma)$.
- ▶ Le **décodeur** est un deep neural network $g(\cdot; \Theta)$.
- ▶ La composition de $f(\cdot; \Phi)$, $s(\cdot; \mu, \Sigma)$ et $g(\cdot; \Theta)$ forme un **réseau de neurones stochastique** $\mathcal{N}(\cdot; \Phi, \Theta)$.



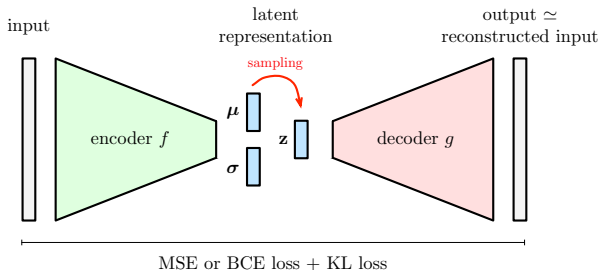
VARIATIONAL AUTOENCODER (VAE)

1. L'encodeur $f(\cdot; \Phi)$ encode une data \mathbf{x} en deux vecteurs $\mu_{\mathbf{x}}$ et $\Sigma_{\mathbf{x}}$ de plus petite dimension.
2. Le sampler $s(\cdot; \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$ sample un point \mathbf{z} dans l'espace latent selon la loi normale $\mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$.
3. Le décodeur $g(\cdot; \Theta)$ reconstruit \mathbf{x} à partir de \mathbf{z} .



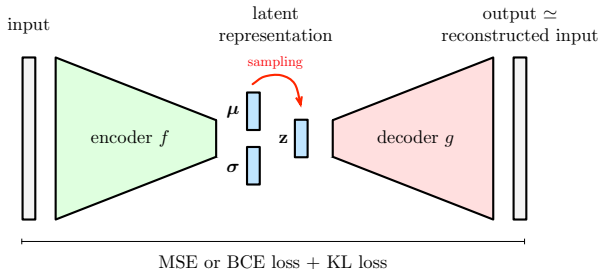
VARIATIONAL AUTOENCODER (VAE)

1. L'**encodeur** $f(\cdot; \Phi)$ encode une data \mathbf{x} en deux vecteurs $\mu_{\mathbf{x}}$ et $\Sigma_{\mathbf{x}}$ de plus petite dimension.
2. Le **sampler** $s(\cdot; \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$ sample un point \mathbf{z} dans l'espace latent selon la loi normale $\mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$.
3. Le **décodeur** $g(\cdot; \Theta)$ reconstruit \mathbf{x} à partir de \mathbf{z} .

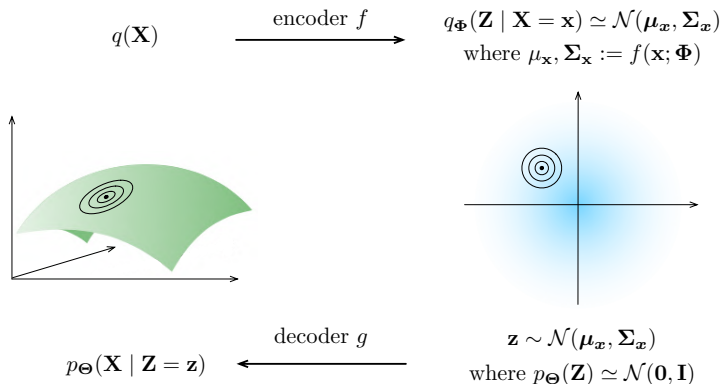


VARIATIONAL AUTOENCODER (VAE)

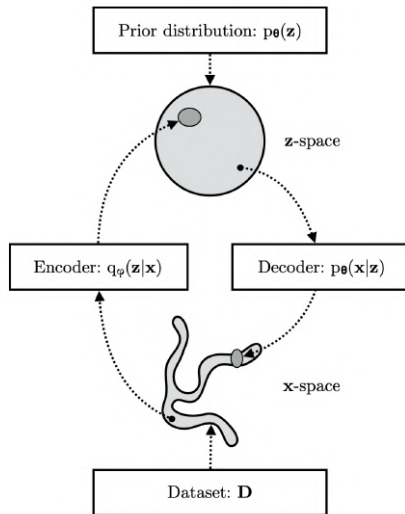
1. L'**encodeur** $f(\cdot; \Phi)$ encode une data \mathbf{x} en deux vecteurs $\mu_{\mathbf{x}}$ et $\Sigma_{\mathbf{x}}$ de plus petite dimension.
2. Le **sampler** $s(\cdot; \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$ sample un point \mathbf{z} dans l'espace latent selon la loi normale $\mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})$.
3. Le **décodeur** $g(\cdot; \Theta)$ reconstruit \mathbf{x} à partir de \mathbf{z} .



VARIATIONAL AUTOENCODER (VAE)



VARIATIONAL AUTOENCODER (VAE)



VARIATIONAL AUTOENCODER (VAE)

- ▶ $q(\mathbf{x})$: Distribution empirique des data.
- ▶ $q_{\Phi}(\mathbf{z} \mid \mathbf{x})$: Distribution des data latentes étant donné la data originale \mathbf{x} . Fonction de \mathbf{z} (variable \mathbf{X} fixé à \mathbf{x}).
→ Dépend de l'encodeur $f(\cdot; \Phi)$.
- ▶ $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$: Distribution des data reconstruites étant donné la data latente \mathbf{z} . Fonction de \mathbf{x} (variable \mathbf{Z} fixé à \mathbf{z}).
→ Dépend du décodeur $g(\cdot; \Theta)$.
- ▶ $p_{\Theta}(\mathbf{x})$: Distribution (marginale) des data reconstruites, induite par les distributions $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$ pour tout \mathbf{z} .
 - La reconstruction s'effectue correctement si $q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ $q(\mathbf{x})$: Distribution empirique des data.
- ▶ $q_{\Phi}(\mathbf{z} \mid \mathbf{x})$: Distribution des data latentes étant donné la data originale \mathbf{x} . Fonction de \mathbf{z} (variable \mathbf{X} fixé à \mathbf{x}).
→ Dépend de l'encodeur $f(\cdot; \Phi)$.
- ▶ $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$: Distribution des data reconstruites étant donné la data latente \mathbf{z} . Fonction de \mathbf{x} (variable \mathbf{Z} fixé à \mathbf{z}).
→ Dépend du décodeur $g(\cdot; \Theta)$.
- ▶ $p_{\Theta}(\mathbf{x})$: Distribution (marginale) des data reconstruites, induite par les distributions $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$ pour tout \mathbf{z} .
- La reconstruction s'effectue correctement si $q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ $q(\mathbf{x})$: Distribution empirique des data.
- ▶ $q_{\Phi}(\mathbf{z} \mid \mathbf{x})$: Distribution des data latentes étant donné la data originale \mathbf{x} . Fonction de \mathbf{z} (variable \mathbf{X} fixé à \mathbf{x}).
→ Dépend de l'encodeur $f(\cdot; \Phi)$.
- ▶ $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$: Distribution des data reconstruites étant donné la data latente \mathbf{z} . Fonction de \mathbf{x} (variable \mathbf{Z} fixé à \mathbf{z}).
→ Dépend du décodeur $g(\cdot; \Theta)$.
- ▶ $p_{\Theta}(\mathbf{x})$: Distribution (marginale) des data reconstruites, induite par les distributions $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$ pour tout \mathbf{z} .
- La reconstruction s'effectue correctement si $q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ $q(\mathbf{x})$: Distribution empirique des data.
- ▶ $q_{\Phi}(\mathbf{z} \mid \mathbf{x})$: Distribution des data latentes étant donné la data originale \mathbf{x} . Fonction de \mathbf{z} (variable \mathbf{X} fixé à \mathbf{x}).
→ Dépend de l'encodeur $f(\cdot; \Phi)$.
- ▶ $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$: Distribution des data reconstruites étant donné la data latente \mathbf{z} . Fonction de \mathbf{x} (variable \mathbf{Z} fixé à \mathbf{z}).
→ Dépend du décodeur $g(\cdot; \Theta)$.
- ▶ $p_{\Theta}(\mathbf{x})$: Distribution (marginale) des data reconstruites, induite par les distributions $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$ pour tout \mathbf{z} .
- La reconstruction s'effectue correctement si $q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$.

VARIATIONAL AUTOENCODER (VAE)

- ▶ $q(\mathbf{x})$: Distribution empirique des data.
- ▶ $q_{\Phi}(\mathbf{z} \mid \mathbf{x})$: Distribution des data latentes étant donné la data originale \mathbf{x} . Fonction de \mathbf{z} (variable \mathbf{X} fixé à \mathbf{x}).
→ Dépend de l'encodeur $f(\cdot; \Phi)$.
- ▶ $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$: Distribution des data reconstruites étant donné la data latente \mathbf{z} . Fonction de \mathbf{x} (variable \mathbf{Z} fixé à \mathbf{z}).
→ Dépend du décodeur $g(\cdot; \Theta)$.
- ▶ $p_{\Theta}(\mathbf{x})$: Distribution (marginale) des data reconstruites, induite par les distributions $p_{\Theta}(\mathbf{x} \mid \mathbf{z})$ pour tout \mathbf{z} .
 - La reconstruction s'effectue correctement si $q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$.

KL-DIVERGENCE ET MAXIMUM LIKELIHOOD

- **Rappel:** étant donné la distribution empirique $q(\mathbf{x})$, on cherche à *apprendre* la distribution $q(\mathbf{x})$, c'est à dire à trouver $p_{\Theta}(\mathbf{x})$ telle que

$$q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$$

- Pour cela, on minimise la *divergence de Kullback–Leibler* entre $q(\mathbf{x})$ et $p_{\Theta}(\mathbf{x})$ définie par:

$$\begin{aligned} D_{\text{KL}}(q \parallel p_{\Theta}) &= \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\Theta}(\mathbf{x})} d\mathbf{x} \\ &= \int q(\mathbf{x}) (\log q(\mathbf{x}) - \log p_{\Theta}(\mathbf{x})) d\mathbf{x} \\ &= \mathbb{E}_{q(\mathbf{x})} [\log q(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] \end{aligned}$$

KL-DIVERGENCE ET MAXIMUM LIKELIHOOD

- **Rappel:** étant donné la distribution empirique $q(\mathbf{x})$, on cherche à *apprendre* la distribution $q(\mathbf{x})$, c'est à dire à trouver $p_{\Theta}(\mathbf{x})$ telle que

$$q(\mathbf{x}) \simeq p_{\Theta}(\mathbf{x})$$

- Pour cela, on minimise la *divergence de Kullback–Leibler* entre $q(\mathbf{x})$ et $p_{\Theta}(\mathbf{x})$ définie par:

$$\begin{aligned} D_{\text{KL}}(q \parallel p_{\Theta}) &= \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p_{\Theta}(\mathbf{x})} d\mathbf{x} \\ &= \int q(\mathbf{x}) (\log q(\mathbf{x}) - \log p_{\Theta}(\mathbf{x})) d\mathbf{x} \\ &= \mathbb{E}_{q(\mathbf{x})} [\log q(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] \end{aligned}$$

VAE: LOSS FUNCTION

- Ceci revient à maximiser la *log-likelihood* de $p_{\Theta}(\mathbf{x})$ définie par:

$$\mathbb{E}_{q(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] = \int q(\mathbf{x}) \log p_{\Theta}(\mathbf{x}) d\mathbf{x}$$

- Dans cette expression, le calcul de $p_{\Theta}(\mathbf{x})$ s'effectue en marginalisant sur la variable latente, i.e.,

$$p_{\Theta}(\mathbf{x}) = \int p_{\Theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\Theta}(\mathbf{x} | \mathbf{z}) p_{\Theta}(\mathbf{z}) d\mathbf{z}$$

- Mais cette expression n'est pas possible à optimiser efficacement (intractable), car son évaluation requiert un parcours de tout l'espace latent pour chaque \mathbf{x} .

VAE: LOSS FUNCTION

- ▶ Ceci revient à maximiser la *log-likelihood* de $p_{\Theta}(\mathbf{x})$ définie par:

$$\mathbb{E}_{q(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] = \int q(\mathbf{x}) \log p_{\Theta}(\mathbf{x}) d\mathbf{x}$$

- ▶ Dans cette expression, le calcul de $p_{\Theta}(\mathbf{x})$ s'effectue en marginalisant sur la variable latente, i.e.,

$$p_{\Theta}(\mathbf{x}) = \int p_{\Theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\Theta}(\mathbf{x} | \mathbf{z}) p_{\Theta}(\mathbf{z}) d\mathbf{z}$$

- ▶ Mais cette expression n'est pas possible à optimiser efficacement (intractable), car son évaluation requiert un parcours de tout l'espace latent pour chaque \mathbf{x} .

VAE: LOSS FUNCTION

- ▶ Ceci revient à maximiser la *log-likelihood* de $p_{\Theta}(\mathbf{x})$ définie par:

$$\mathbb{E}_{q(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] = \int q(\mathbf{x}) \log p_{\Theta}(\mathbf{x}) d\mathbf{x}$$

- ▶ Dans cette expression, le calcul de $p_{\Theta}(\mathbf{x})$ s'effectue en marginalisant sur la variable latente, i.e.,

$$p_{\Theta}(\mathbf{x}) = \int p_{\Theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\Theta}(\mathbf{x} | \mathbf{z}) p_{\Theta}(\mathbf{z}) d\mathbf{z}$$

- ▶ Mais cette expression n'est pas possible à optimiser efficacement (intractable), car son évaluation requiert un parcours de tout l'espace latent pour chaque \mathbf{x} .

VAE: LOSS FUNCTION

- ▶ Par contre, on peut montrer que (cf. notes manuscrites):

$$\begin{aligned} \log p_{\Theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \underbrace{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}_{\sim \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})} \quad \underbrace{\phantom{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \end{aligned} \quad (1)$$

Le terme de droite se nomme *evidence lower bound (ELBO)*.

- ▶ En prenant l'espérance sur $q_{\Phi}(\mathbf{x})$ à gauche et à droite, on a:

$$\begin{aligned} \mathbb{E}_{q_{\Phi}(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}, \mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \mathbb{E}_{q_{\Phi}(\mathbf{x})} \left[D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right) \right] \end{aligned} \quad (2)$$

- ▶ Ainsi, maximiser les termes de droite dans (1) et (2) – ou alors minimiser leurs opposés – permet de maximiser les log-probabilité et log-likelihood de gauche, respectivement.

VAE: LOSS FUNCTION

- ▶ Par contre, on peut montrer que (cf. notes manuscrites):

$$\begin{aligned}\log p_{\Theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \text{D}_{\text{KL}} \left(\underbrace{q_{\Phi}(\mathbf{z} | \mathbf{x})}_{\sim \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})} \parallel \underbrace{p_{\Theta}(\mathbf{z})}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \right)\end{aligned}\quad (1)$$

Le terme de droite se nomme *evidence lower bound (ELBO)*.

- ▶ En prenant l'espérance sur $q_{\Phi}(\mathbf{x})$ à gauche et à droite, on a:

$$\begin{aligned}\mathbb{E}_{q_{\Phi}(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}, \mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \mathbb{E}_{q_{\Phi}(\mathbf{x})} \left[\text{D}_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right) \right]\end{aligned}\quad (2)$$

- ▶ Ainsi, maximiser les termes de droite dans (1) et (2) – ou alors minimiser leurs opposés – permet de maximiser les log-probabilité et log-likelihood de gauche, respectivement.

VAE: LOSS FUNCTION

- ▶ Par contre, on peut montrer que (cf. notes manuscrites):

$$\begin{aligned}\log p_{\Theta}(\mathbf{x}) &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \underbrace{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}_{\sim \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}})} \quad \underbrace{\phantom{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}}_{\sim \mathcal{N}(\mathbf{0}, \mathbf{I})}\end{aligned}\quad (1)$$

Le terme de droite se nomme *evidence lower bound (ELBO)*.

- ▶ En prenant l'espérance sur $q_{\Phi}(\mathbf{x})$ à gauche et à droite, on a:

$$\begin{aligned}\mathbb{E}_{q_{\Phi}(\mathbf{x})} [\log p_{\Theta}(\mathbf{x})] &\geq \mathbb{E}_{q_{\Phi}(\mathbf{z}, \mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right] \\ &- \mathbb{E}_{q_{\Phi}(\mathbf{x})} \left[D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right) \right]\end{aligned}\quad (2)$$

- ▶ Ainsi, maximiser les termes de droite dans (1) et (2) – ou alors minimiser leurs opposés – permet de maximiser les log-probabilité et log-likelihood de gauche, respectivement.

VAE: LOSS FUNCTION

- On prend alors comme fonction de loss individuelle (cf. Eq. (1)):

$$\ell_i(\Phi, \Theta) = \underbrace{-\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right]}_{\text{reconstruction term}} + \underbrace{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}_{\text{regularization term}}$$

- Et comme fonction de loss collective (cf. Eq. (2)):

$$\mathcal{L}(\Phi, \Theta) = \underbrace{-\mathbb{E}_{q_{\Phi}(\mathbf{z}, \mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right]}_{\text{reconstruction term}} + \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{x})} \left[D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right) \right]}_{\text{regularization term}}$$

VAE: LOSS FUNCTION

- On prend alors comme fonction de loss individuelle (cf. Eq. (1)):

$$\ell_i(\Phi, \Theta) = \underbrace{-\mathbb{E}_{q_{\Phi}(\mathbf{z}|\mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right]}_{\text{reconstruction term}} + \underbrace{D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right)}_{\text{regularization term}}$$

- Et comme fonction de loss collective (cf. Eq. (2)):

$$\mathcal{L}(\Phi, \Theta) = \underbrace{-\mathbb{E}_{q_{\Phi}(\mathbf{z}, \mathbf{x})} \left[\log p_{\Theta}(\mathbf{x} | \mathbf{z}) \right]}_{\text{reconstruction term}} + \underbrace{\mathbb{E}_{q_{\Phi}(\mathbf{x})} \left[D_{\text{KL}} \left(q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p_{\Theta}(\mathbf{z}) \right) \right]}_{\text{regularization term}}$$

VAE: LOSS FUNCTION

- Pour un batch de data $\{\mathbf{x}_i : i = 1, \dots, B\}$, on a les *estimateurs* suivants de ces fonctions de loss:

$$\hat{\ell}_i(\Phi, \Theta) = -\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i) + \hat{D}_{\text{KL}}(q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) || p_{\Theta}(\mathbf{z}_i))$$

$$\hat{\mathcal{L}}(\Phi, \Theta) = \sum_{i=1}^B \left(-\log p_{\Theta}(\mathbf{x}_i | \mathbf{z}_i) + \hat{D}_{\text{KL}}(q_{\Phi}(\mathbf{z}_i | \mathbf{x}_i) || p_{\Theta}(\mathbf{z}_i)) \right)$$

où $\mathbf{z}_i \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}_i}, \boldsymbol{\Sigma}_{\mathbf{x}_i})$ avec $\boldsymbol{\mu}_{\mathbf{x}_i}, \boldsymbol{\Sigma}_{\mathbf{x}_i} := f(\mathbf{x}_i; \Phi)$

- Cette loss (collective) est minimisée par backpropagation.

VARIATIONAL AUTOENCODER (VAE): TRAINING

- ▶ Mais le processus de sampling n'est pas différentiable, ce qui pose problème pour l'algorithme de backpropagation.
- ▶ Pour pallier cela, on recourt au “reparametrization trick” et au “log-var trick”.
- ▶ Au lieu de sampler de \mathbf{z} comme suit:

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} \sim q_{\Phi}(\mathbf{Z} \mid \mathbf{X} = \mathbf{x}) = \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}}) \end{cases}$$

on procède de la manière suivante:

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} = \mu_{\mathbf{x}} + \sigma_{\mathbf{x}} \odot \epsilon \text{ où } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ et } \sigma_{\mathbf{x}} = \text{diag}(\Sigma_{\mathbf{x}})^{\frac{1}{2}} \end{cases}$$

VARIATIONAL AUTOENCODER (VAE): TRAINING

- ▶ Mais le processus de sampling n'est pas différentiable, ce qui pose problème pour l'algorithme de backpropagation.
- ▶ Pour pallier cela, on recourt au “reparametrization trick” et au “log-var trick”.
- ▶ Au lieu de sampler de \mathbf{z} comme suit:

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} \sim q_{\Phi}(\mathbf{Z} \mid \mathbf{X} = \mathbf{x}) = \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}}) \end{cases}$$

on procède de la manière suivante:

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} = \mu_{\mathbf{x}} + \sigma_{\mathbf{x}} \odot \epsilon \text{ où } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ et } \sigma_{\mathbf{x}} = \text{diag}(\Sigma_{\mathbf{x}})^{\frac{1}{2}} \end{cases}$$

VARIATIONAL AUTOENCODER (VAE): TRAINING

- ▶ Mais le processus de sampling n'est pas différentiable, ce qui pose problème pour l'algorithme de backpropagation.
- ▶ Pour pallier cela, on recourt au “reparametrization trick” et au “log-var trick”.
- ▶ Au lieu de sampler de \mathbf{z} comme suit:

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} \sim q_{\Phi}(\mathbf{Z} \mid \mathbf{X} = \mathbf{x}) = \mathcal{N}(\mu_{\mathbf{x}}, \Sigma_{\mathbf{x}}) \end{cases}$$

on procède de la manière suivante:

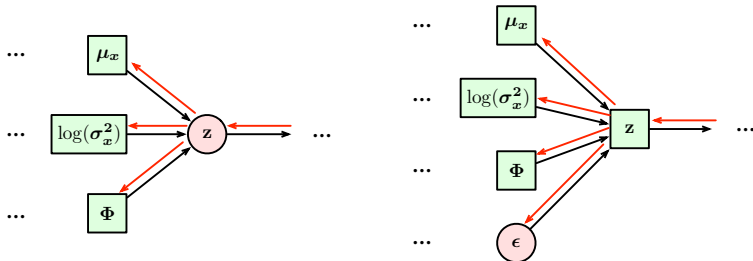
$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} = \mu_{\mathbf{x}} + \sigma_{\mathbf{x}} \odot \epsilon \text{ où } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \text{ et } \sigma_{\mathbf{x}} = \text{diag}(\Sigma_{\mathbf{x}})^{\frac{1}{2}} \end{cases}$$

VARIATIONAL AUTOENCODER (VAE): TRAINING

- ▶ En pratique, pour des raisons de stabilité, on apprend les vecteurs $\mu_{\mathbf{x}}$ et $\log(\sigma_{\mathbf{x}}^2)$ au lieu de $\mu_{\mathbf{x}}$ et $\sigma_{\mathbf{x}}^2$. On a alors

$$\begin{cases} \mu_{\mathbf{x}}, \Sigma_{\mathbf{x}} := f(\mathbf{x}; \Phi) \\ \mathbf{z} = \mu_{\mathbf{x}} + \sigma_{\mathbf{x}} \odot \epsilon = \mu_{\mathbf{x}} + e^{\frac{1}{2}\log(\sigma_{\mathbf{x}}^2)} \odot \epsilon \quad \text{où } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{cases}$$

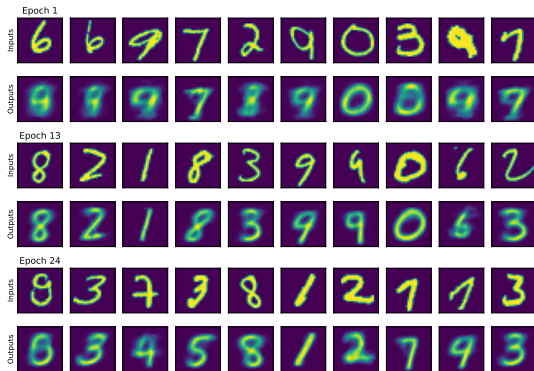
VARIATIONAL AUTOENCODER (VAE): TRAINING



Computational graphs with and without the reparametrization trick. Green and red nodes are deterministic and stochastic nodes, respectively. Red arrows illustrate backpropagation.

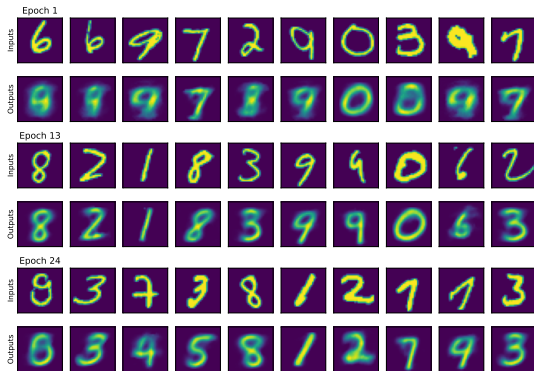
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ L'espace latent est réparti de manière bien plus gaussienne et les capacité de génération sont censées être améliorées (poids des termes de reconstruction et régularisation: 10 et 0.001).
- ▶ Pas si clair sur un espace latent de dim 2 (cf. slides suivants)...



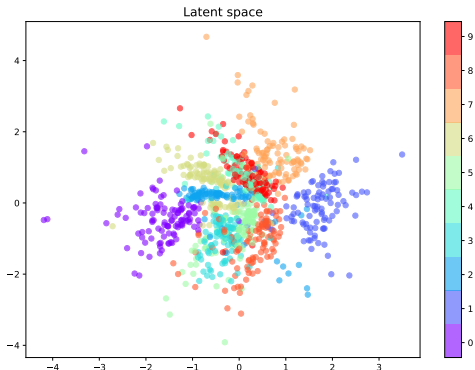
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- L'espace latent est réparti de manière bien plus gaussienne et les capacité de génération sont censées être améliorées (poids des termes de reconstruction et régularisation: 10 et 0.001).
- Pas si clair sur un espace latent de dim 2 (cf. slides suivants)...



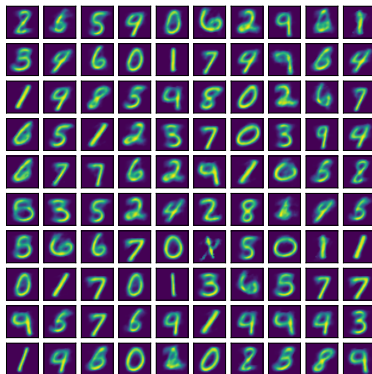
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ L'espace latent est réparti de manière bien plus gaussienne et les capacité de génération sont censées être améliorées (poids des termes de reconstruction et régularisation: 10 et 0.001).
- ▶ Pas si clair sur un espace latent de dim 2 (cf. slides suivants)...



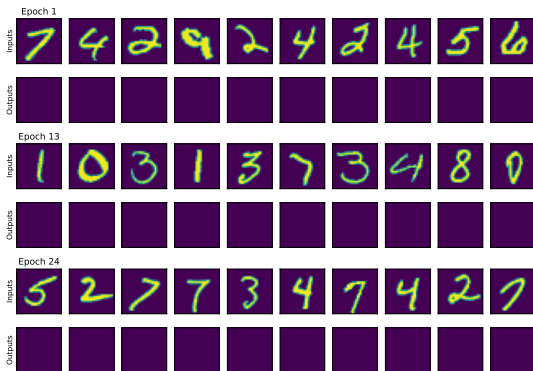
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ L'espace latent est réparti de manière bien plus gaussienne et les capacité de génération sont censées être améliorées (poids des termes de reconstruction et régularisation: 10 et 0.001).
- ▶ Pas si clair sur un espace latent de dim 2 (cf. slides suivants)...



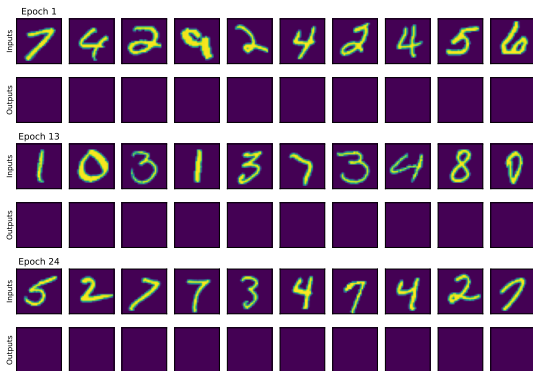
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si on annule le terme de reconstruction, les points de l'espace latent sont distribués de manière parfaitement normale...
- ▶ mais le réseau n'a rien appris à reconstruire.



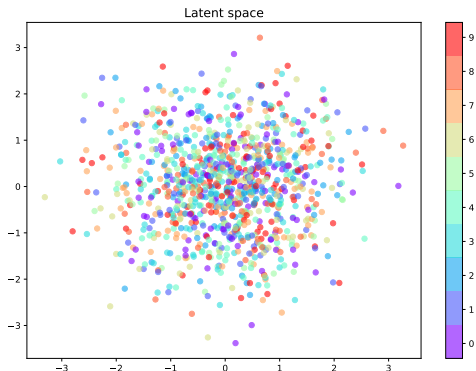
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si on annule le terme de reconstruction, les points de l'espace latent sont distribués de manière parfaitement normale...
- ▶ mais le réseau n'a rien appris à reconstruire.



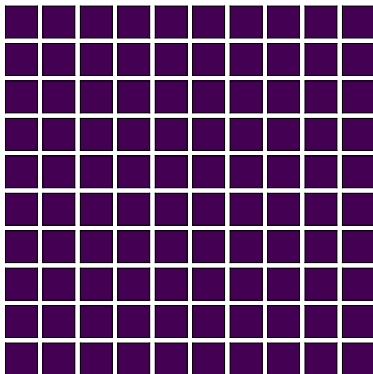
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si on annule le terme de reconstruction, les points de l'espace latent sont distribués de manière parfaitement normale...
- ▶ mais le réseau n'a rien appris à reconstruire.



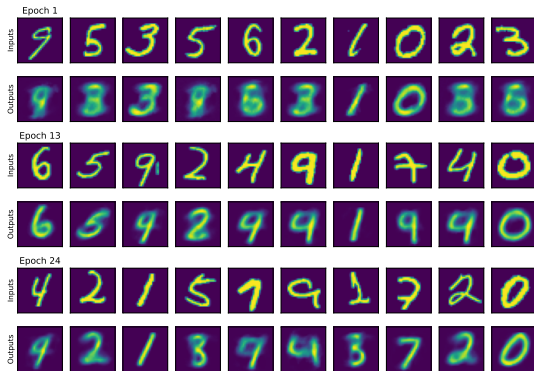
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si on annule le terme de reconstruction, les points de l'espace latent sont distribués de manière parfaitement normale...
- ▶ mais le réseau n'a rien appris à reconstruire.



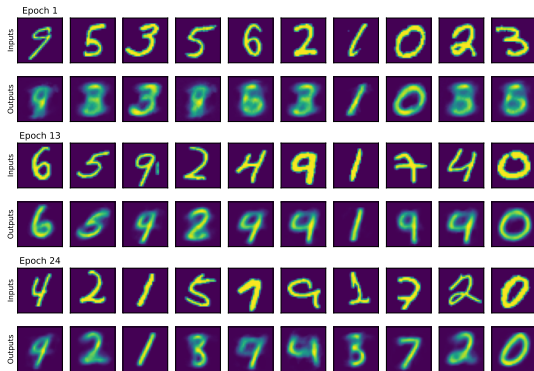
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si, au contraire, on annule le terme de régularisation, les points de l'espace latent ne sont plus distribués de manière normale. . .
- ▶ le réseau a appris à reconstruire des data, mais en samplant dans $\mathcal{N}(0, 1)$, on génère des data assez similaires.



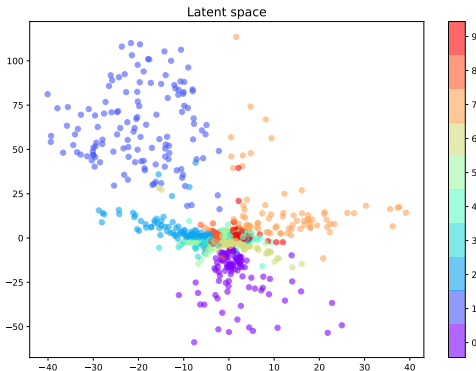
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si, au contraire, on annule le terme de régularisation, les points de l'espace latent ne sont plus distribués de manière normale. . .
- ▶ le réseau a appris à reconstruire des data, mais en samplant dans $\mathcal{N}(0, 1)$, on génère des data assez similaires.



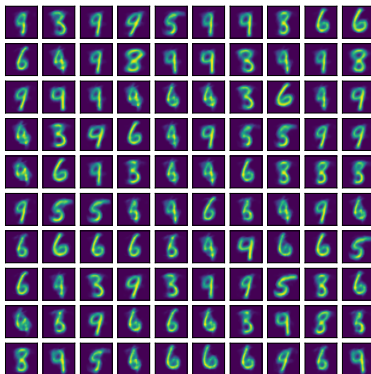
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si, au contraire, on annule le terme de régularisation, les points de l'espace latent ne sont plus distribués de manière normale. . .
- ▶ le réseau a appris à reconstruire des data, mais en samplant dans $\mathcal{N}(0, 1)$, on génère des data assez similaires.



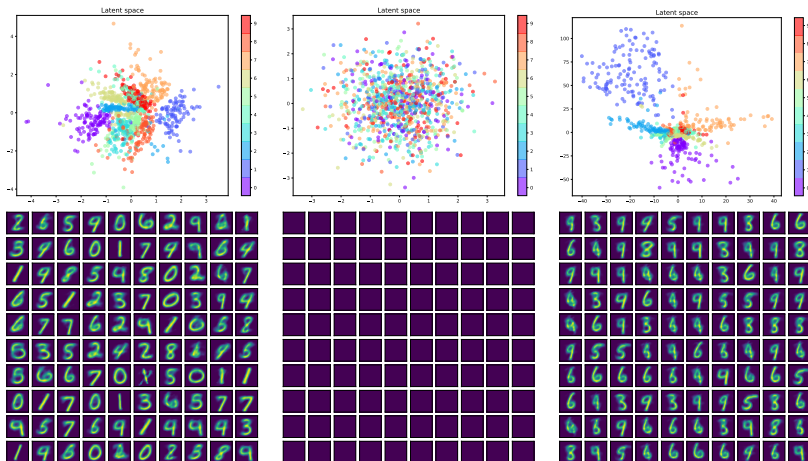
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

- ▶ Si, au contraire, on annule le terme de régularisation, les points de l'espace latent ne sont plus distribués de manière normale. . .
- ▶ le réseau a appris à reconstruire des data, mais en samplant dans $\mathcal{N}(0, 1)$, on génère des data assez similaires.



VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 2

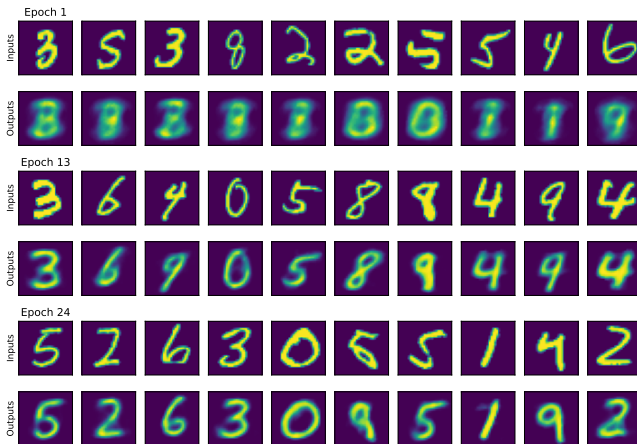
- ▶ Comparaison entre les trois situations.



Left: combination of MSE and KL losses; **Center:** only KL loss; **Right:** only MSE loss.

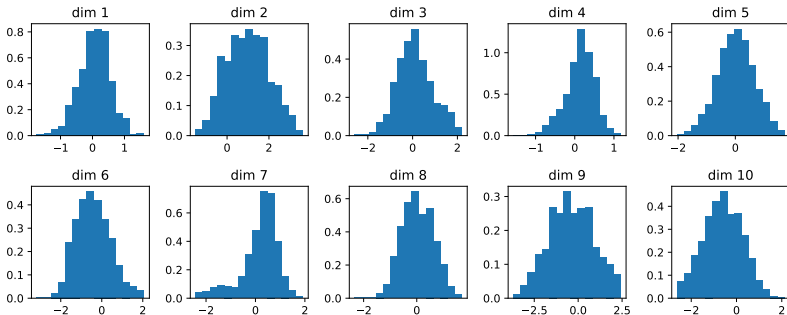
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 10

- En augmentant la dimension de l'espace latent, la génération de data est de meilleure qualité.



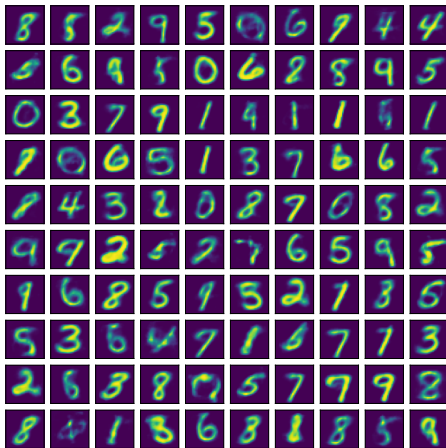
VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 10

- En augmentant la dimension de l'espace latent, la génération de data est de meilleure qualité.

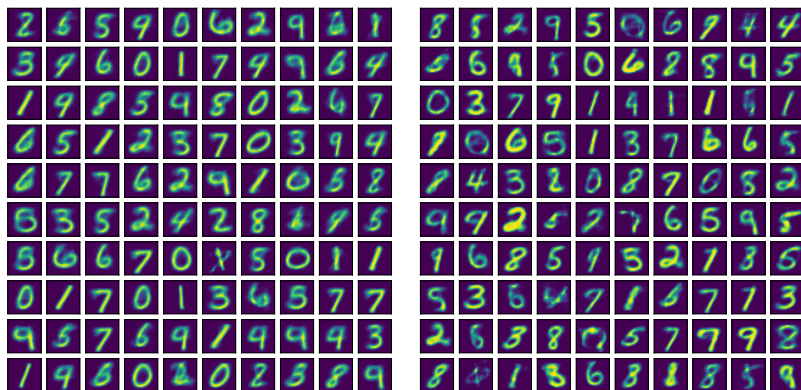


VAE LINÉAIRE: ESPACE LATENT DE DIMENSION 10

- En augmentant la dimension de l'espace latent, la génération de data est de meilleure qualité.



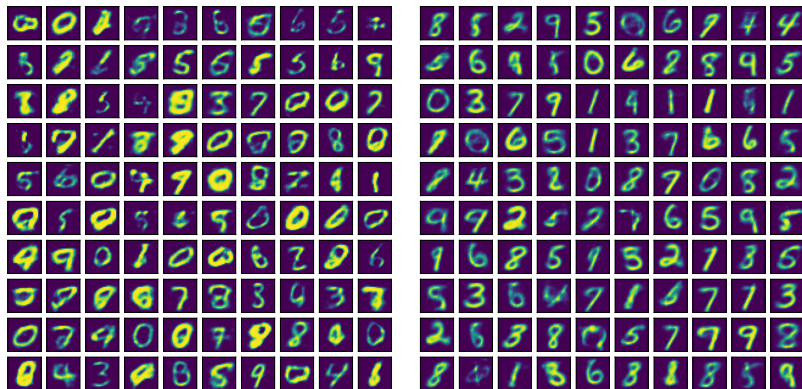
VAE LINÉAIRE: ESPACES LATENTS DE DIMENSIONS 2 ET 10



Data generation of two VAEs with latent spaces of dimensions 2 (left) and 10 (right), respectively.

AE vs VAE: ESPACE LATENT DE DIMENSION 10

► Génération de data: AE vs VAE.



Data generation of an AE (left) and a VAE (right) with latent space of dimension 10.

CONCLUSION

- ▶ Les autoencodeurs variationnels (VAEs) améliorent les capacité de génération de data par rapports aux autoencodeurs (AEs).
- ▶ L'ajout de stochasticité et l'imposition d'une distribution sur l'espace latent répond à des contraintes de continuité et de complétude sur ce espace.
- ▶ La continuité de l'espace latent rend les variables latentes interprétables: on comprend bien l'effet induit par la modification d'une variable alors que les autres restent fixes.

CONCLUSION

- ▶ Les autoencodeurs variationnels (VAEs) améliorent les capacité de génération de data par rapports aux autoencodeurs (AEs).
- ▶ L'ajout de stochasticité et l'imposition d'une distribution sur l'espace latent répond à des contraintes de continuité et de complétude sur ce espace.
- ▶ La continuité de l'espace latent rend les variables latentes interprétables: on comprend bien l'effet induit par la modification d'une variable alors que les autres restent fixes.

CONCLUSION

- ▶ Les autoencodeurs variationnels (VAEs) améliorent les capacité de génération de data par rapports aux autoencodeurs (AEs).
- ▶ L'ajout de stochasticité et l'imposition d'une distribution sur l'espace latent répond à des contraintes de continuité et de complétude sur ce espace.
- ▶ La continuité de l'espace latent rend les variables latentes interprétables: on comprend bien l'effet induit par la modification d'une variable alors que les autres restent fixes.

BIBLIOGRAPHIE



Alexander Amini (2022).
Alexander amini: Youtube channel.



CNRS - Formation FIDLE (2022).
Cnrs - formation fidle: Youtube channel.



Fleuret, F. (2022).
Deep Learning Course.



Kingma, D. P. and Welling, M. (2014).
Auto-encoding variational bayes.
In Bengio, Y. and LeCun, Y., editors, *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.



Kingma, D. P. and Welling, M. (2019).
An introduction to variational autoencoders.
Found. Trends Mach. Learn., 12(4):307–392.



Sebastian Raschka (2022).
Sebastian raschka: Youtube channel.



Stephen Odaibo (2020).
Stephen odaibo : Meduim post.