

Modèles génératifs profonds

Guillaume Bourmaud

PLAN

I. Introduction

II. Réseau inversible (NF)

III. Auto-encodeur variationnel (VAE)

IV. Réseau antagoniste (GAN)

I) Introduction

Principe d'un modèle génératif profond

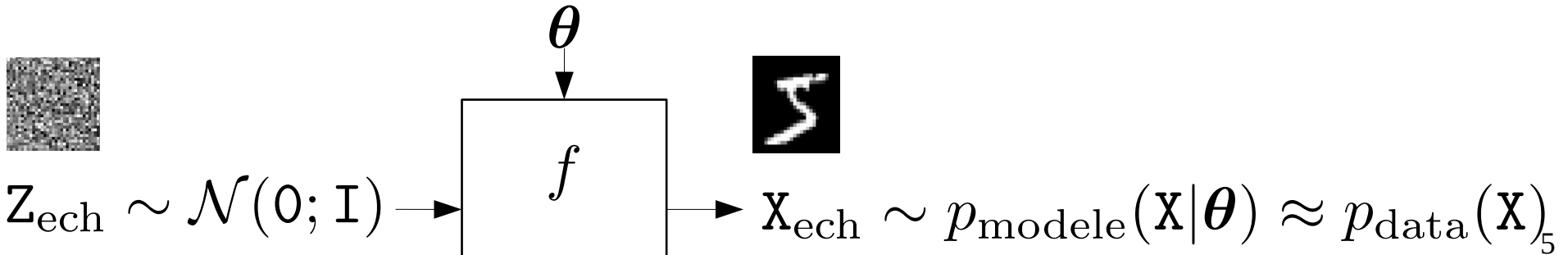
Base de données **non-étiquetées** : $\{X_i\}_{i=1\dots N}$

“est un échantillon de”

Point de vue probabiliste : $X_i \sim p_{\text{data}}(X)$ **Inconnue**



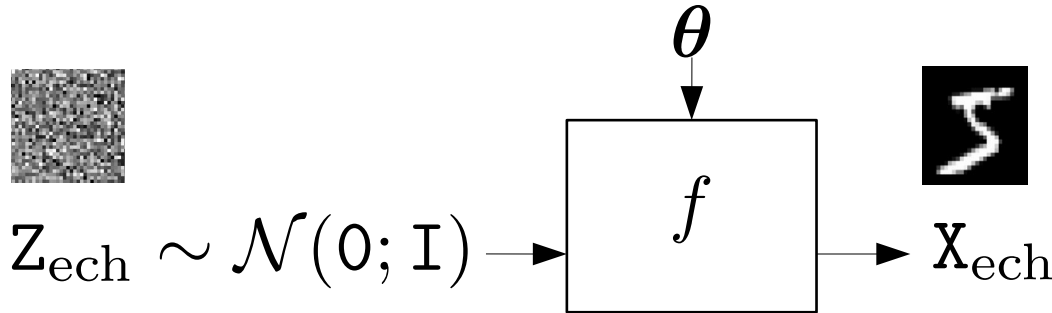
Objectif : Optimiser les paramètres d'un réseau de neurones profond afin de générer de **nouveaux échantillons** de $p_{\text{data}}(X)$



1)

Difficulté de la tâche d'apprentissage

Base de données **non-étiquetées** : $\{X_i\}_{i=1\dots N}$

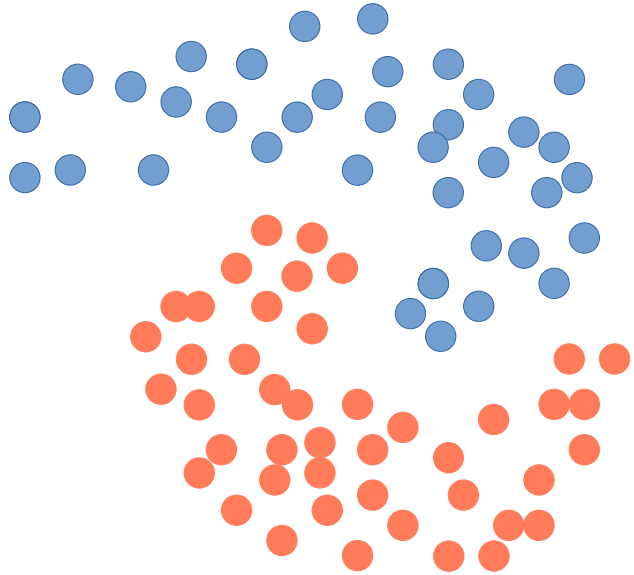


On ne dispose pas de couples $\{Z_i, X_i\}_{i=1\dots N}$

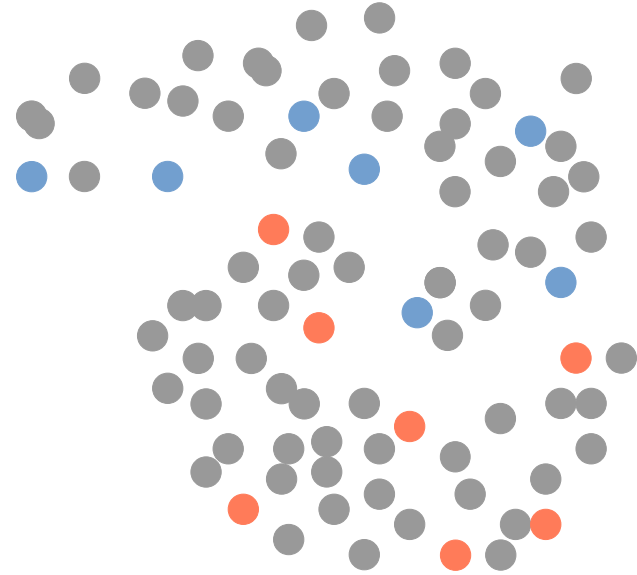
→ problème d'apprentissage non-supervisé

1)

Exemple d'utilisation : Apprentissage semi-supervisé



Supervisé



Semi-Supervisé

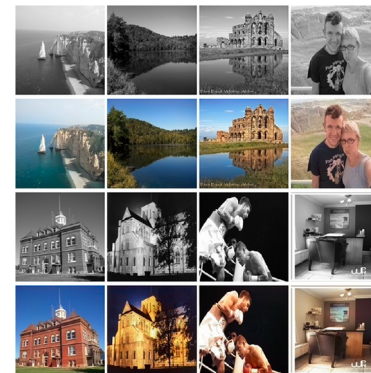
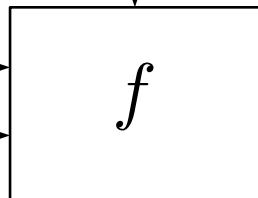
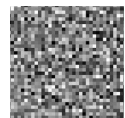
1)

Principe d'un modèle génératif profond conditionnel

Base de données **étiquetées** : $\{Y_i, X_i\}_{i=1\dots N}$

Point de vue probabiliste : $X_i \sim p_{\text{data}}(X|Y_i)$ **Inconnue**

Objectif : Optimiser les paramètres d'un réseau de neurones afin de générer, à partir d'un Y , de **nouveaux échantillons** de $p_{\text{data}}(X|Y)$


 θ
 Y
 $Z_{\text{ech}} \sim \mathcal{N}(0; I)$

 $X_{\text{ech}} \sim p_{\text{modele}}(X|Y, \theta) \approx p_{\text{data}}(X|Y)$

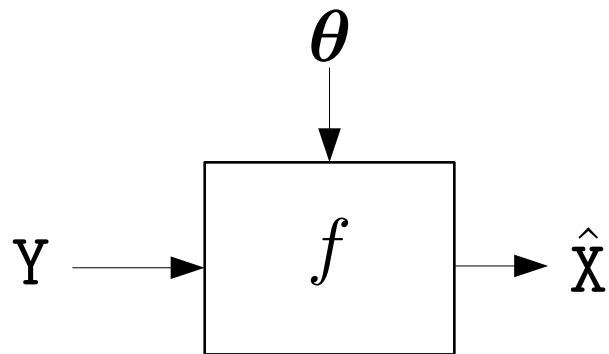
1)

Différence vis-à-vis de l'apprentissage supervisé



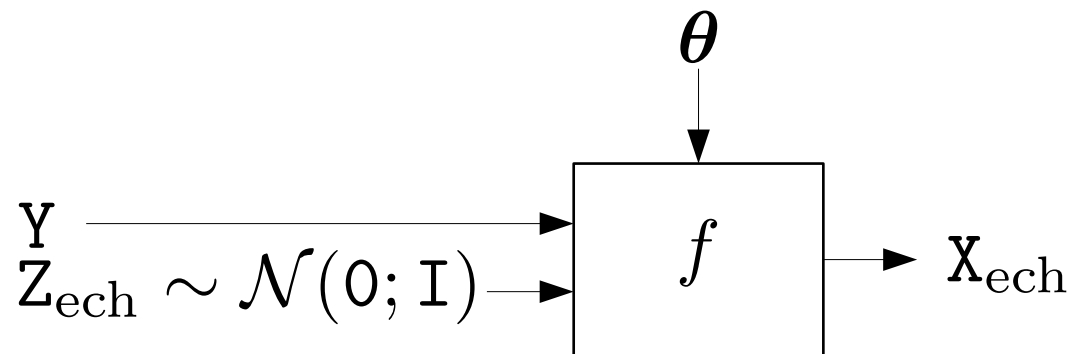
$$p_{\text{data}}(\mathbf{X}|\mathbf{Y})$$

de forme complexe
(ex : multimodale)



Supervisé → Régresseur ou classifieur

$p_{\text{data}}(\mathbf{X}|\mathbf{Y})$ de forme "simple"
(ex : unimodale)



Modèle génératif → Échantillonneur

$p_{\text{data}}(\mathbf{X}|\mathbf{Y})$ de forme "complexe"

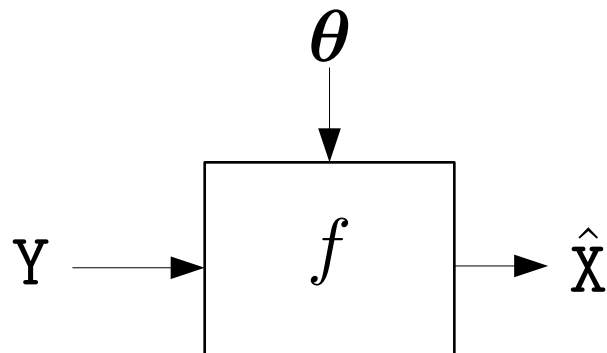
I)

Différence vis-à-vis de l'apprentissage supervisé (suite)



$$p_{\text{data}}(\mathbf{X}|\mathbf{Y})$$

de forme complexe
(ex : multimodale)



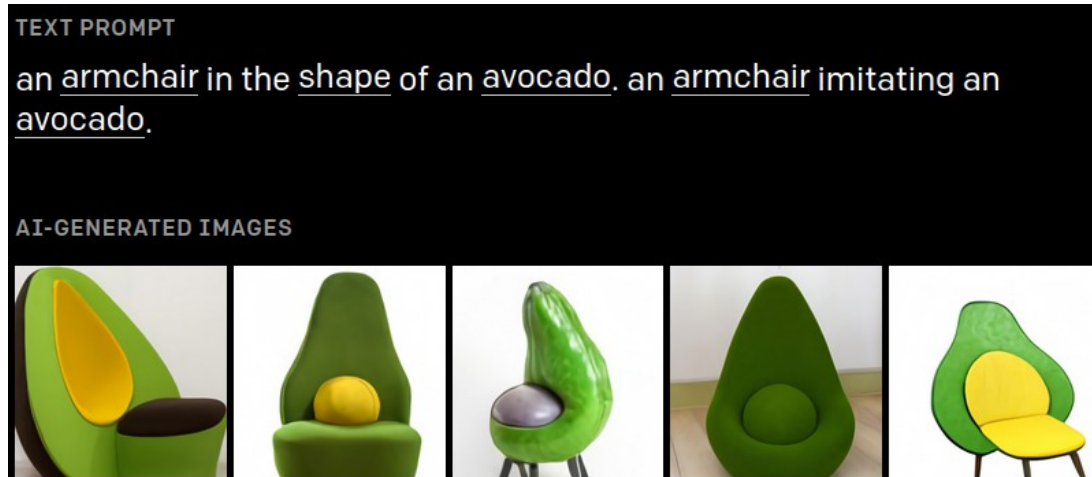
**Question : Que se passe-t-il si on apprend
un régresseur pour colorer une image ?**

Supervisé → Régresseur ou classifieur

$p_{\text{data}}(\mathbf{X}|\mathbf{Y})$ de forme "simple"
(ex : unimodale)

Exemples d'utilisation

- Colorisation d'image
- Super-résolution
- Édition/Synthèse de données (ex : retouche d'image)



<https://openai.com/blog/dall-e/>

Génération de faux contenus :

- faux article de presse
- fausse vidéo
- faux signal sonore

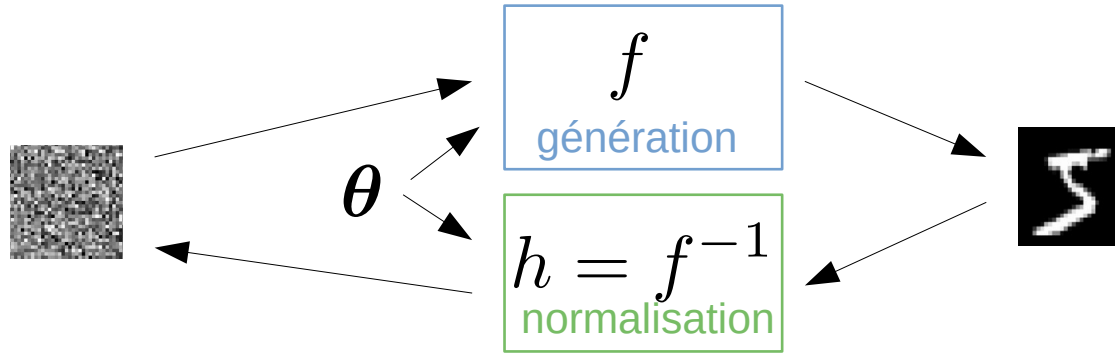
Comment entraîner un modèle génératif profond ?

Beaucoup de solutions possibles, nous allons en étudier trois :

- le réseau inversible (NF)
- l'auto-encodeur variationnel (VAE)
- le réseau antagoniste (GAN)

II) Réseau inversible (NF)

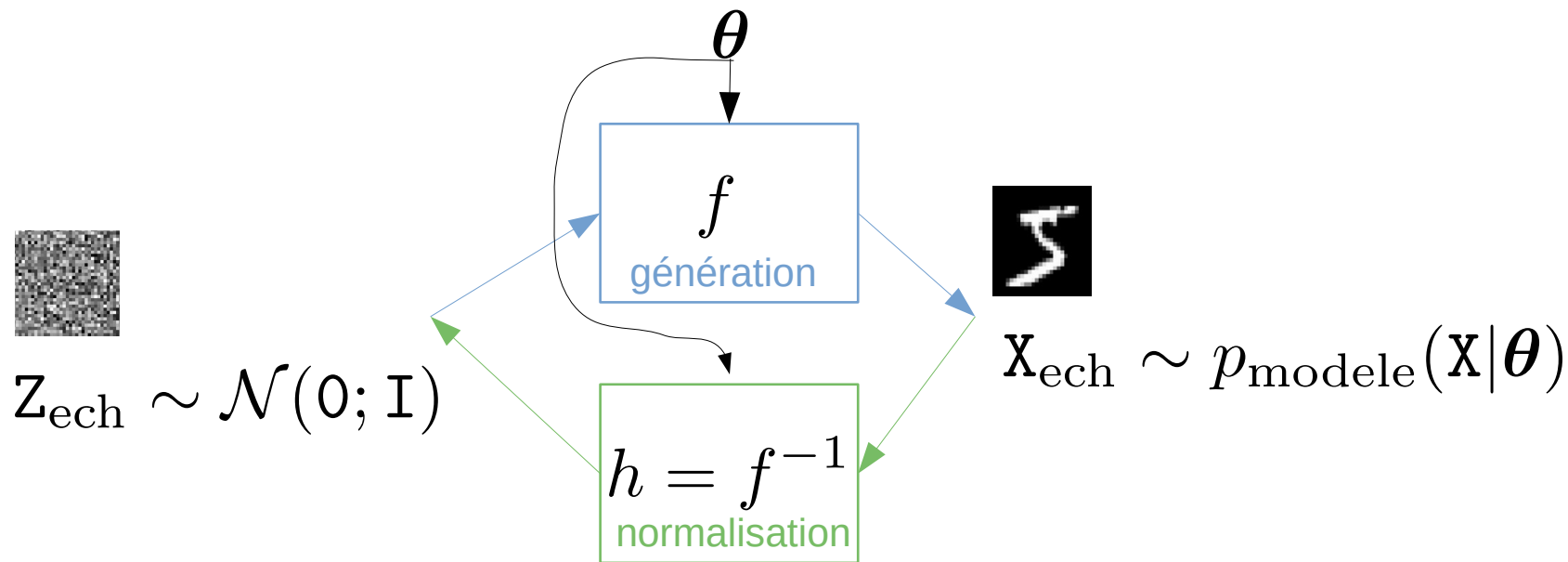
Réseau inversible : idée générale



Ne pas disposer de couples $\{Z_i, X_i\}_{i=1\dots N}$ n'est plus un problème.

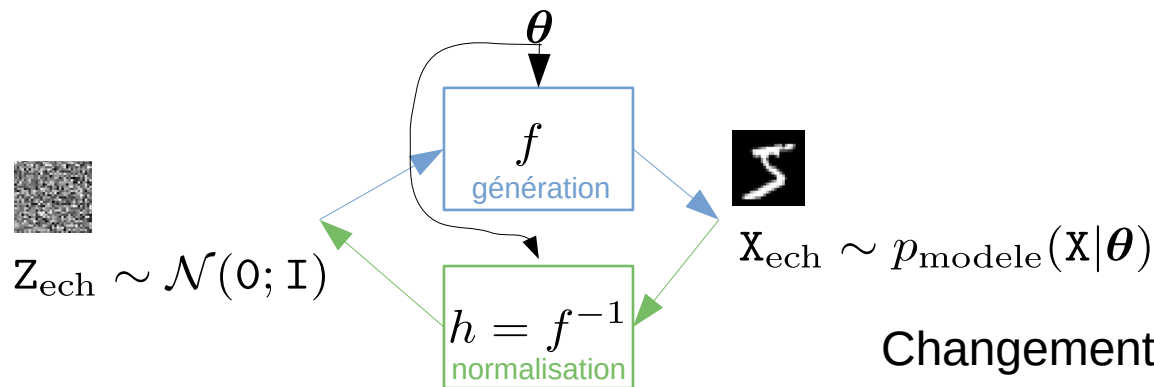
II)

Réseau inversible : idée générale (suite)



II)

Réseau inversible : formalisme



Changement de variable : $Z = h(X; \theta)$

$$1 = \int p(Z) dZ = \int \mathcal{N}(Z; 0, \mathbf{I}) dZ = \int \underbrace{\mathcal{N}(h(X; \theta); 0, \mathbf{I}) |\det J_h(X)|}_{p_{modele}(X|\theta)} dX$$

$$-\ln p_{modele}(X|\theta) = ||h(X; \theta)||^2 - \ln |\det J_h(X; \theta)| + \text{cst}_{\theta}$$

Réseau inversible : apprentissage

Objectif : optimiser θ pour que $p_{\text{modele}}(\mathbf{X}|\theta) \approx p_{\text{data}}(\mathbf{X})$

Nécessité de choisir une fonction de similarité (“divergence”).

Pour un NF on choisit généralement la divergence de Kullback-Leibler suivante :

$$KL(p_{\text{data}}(\mathbf{X})||p_{\text{modele}}(\mathbf{X}|\theta)) = - \int p_{\text{data}}(\mathbf{X}) \ln p_{\text{modele}}(\mathbf{X}|\theta) d\mathbf{X} + \text{cst}_{\theta}$$

$$L(\theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x})}(-\ln p_{\text{modele}}(\mathbf{X}|\theta)) \approx \frac{1}{N} \sum_{i=1}^N -\ln p_{\text{modele}}(\mathbf{x}_i|\theta)$$

= maximisation de la vraisemblance

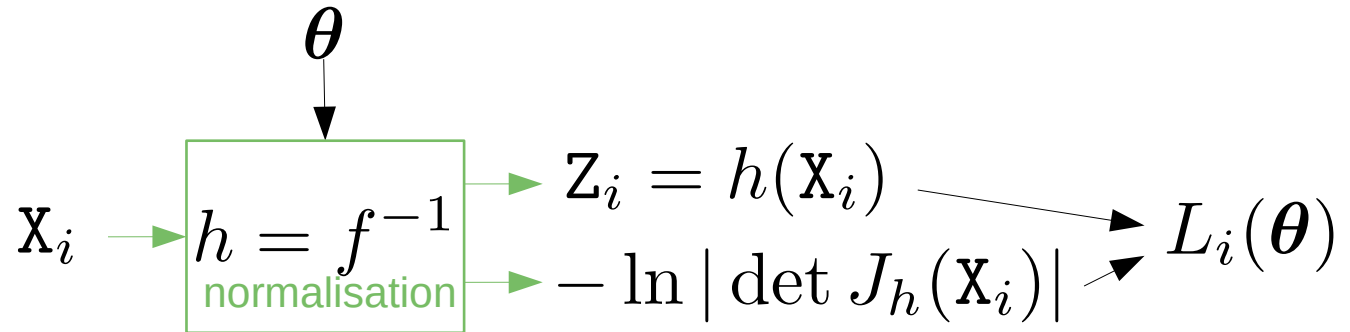
Réseau inversible : apprentissage (suite)

$$L(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N ||h(\mathbf{x}_i; \boldsymbol{\theta})||^2 - \ln |\det J_h(\mathbf{x}_i; \boldsymbol{\theta})|$$

Essaye de transformer \mathbf{X}_i
proche du tenseur "zéro"

Empêche que tout le monde se
transforme en un tenseur "zéro"

Descente de gradient



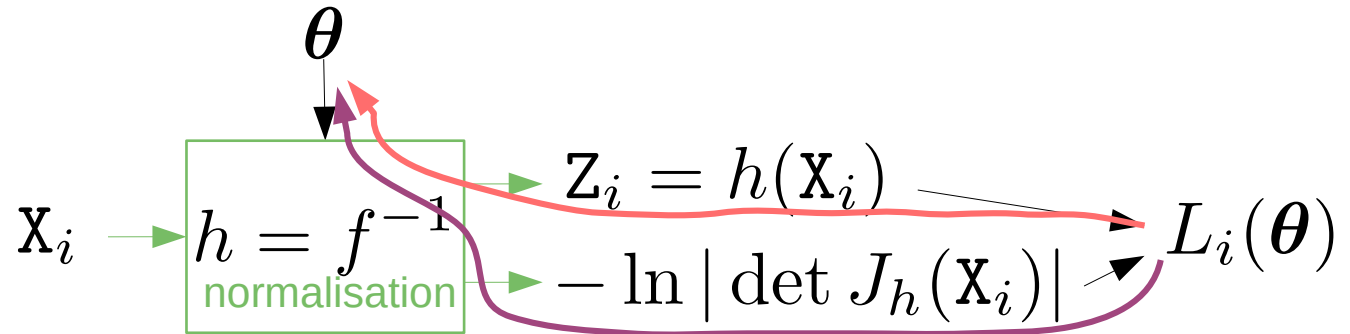
Réseau inversible : apprentissage (suite)

$$L(\boldsymbol{\theta}) \approx \frac{1}{N} \sum_{i=1}^N ||h(\mathbf{x}_i; \boldsymbol{\theta})||^2 - \ln |\det J_h(\mathbf{x}_i; \boldsymbol{\theta})|$$

Essaye de transformer \mathbf{X}_i
proche du tenseur "zéro"

Empêche que tout le monde se
transforme en un tenseur "zéro"

Descente de gradient



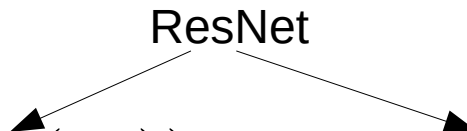
En pratique : arrêt prématuré ("early stopping") sur un ensemble de validation.

Réseau inversible : architecture

Besoin d'une architecture très spécifique, chaque couche doit :

1. avoir la même taille en entrée et en sortie
2. être inversible, et pour pouvoir échantillonner efficacement cette fonction inverse doit se calculer efficacement
3. avoir comme propriété que le “logdet” de sa jacobienne se calcule efficacement et soit différentiable.

Couche la plus souvent utilisée : “Affine coupling layer”

$$\begin{array}{l} X_1 = Z_1 \\ X_2 = \exp(s_{\theta}(Z_1)) \odot Z_2 + m_{\theta}(Z_1) \end{array}$$


jacobienne triangulaire !

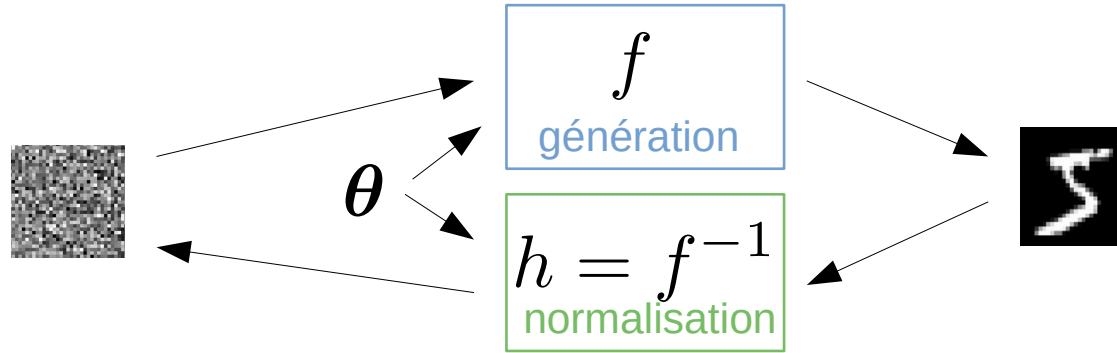
Réseau inversible : avantages et inconvénients

- + capable d'échantillonner efficacement
- + capable de calculer la (log-)probabilité d'une donnée
- + apprentissage possible par maximum de vraisemblance
- contrainte d'architecture inversible
- contrainte du même nombre de dimension entrée/sortie
- la forme de la distribution qu'on transforme est fortement contrainte
- contrainte sur le calcul du log-déterminant

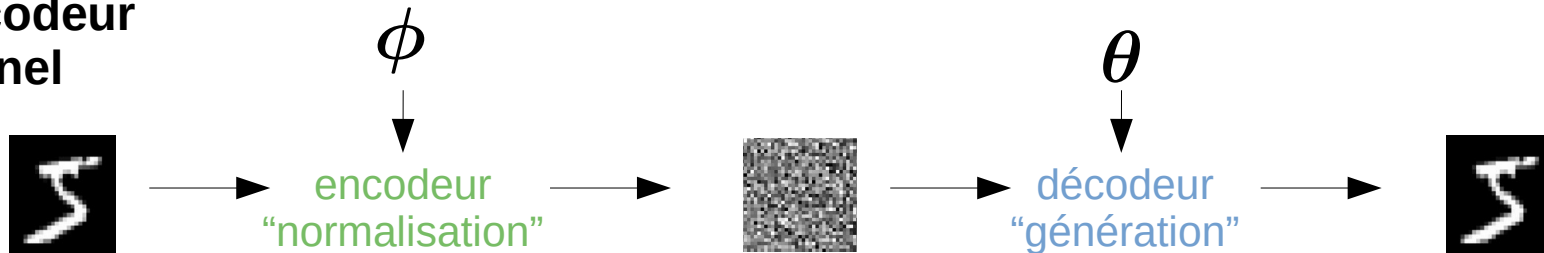
III) Auto-encodeur variationnel (VAE)

Auto-encodeur variationnel : idée générale

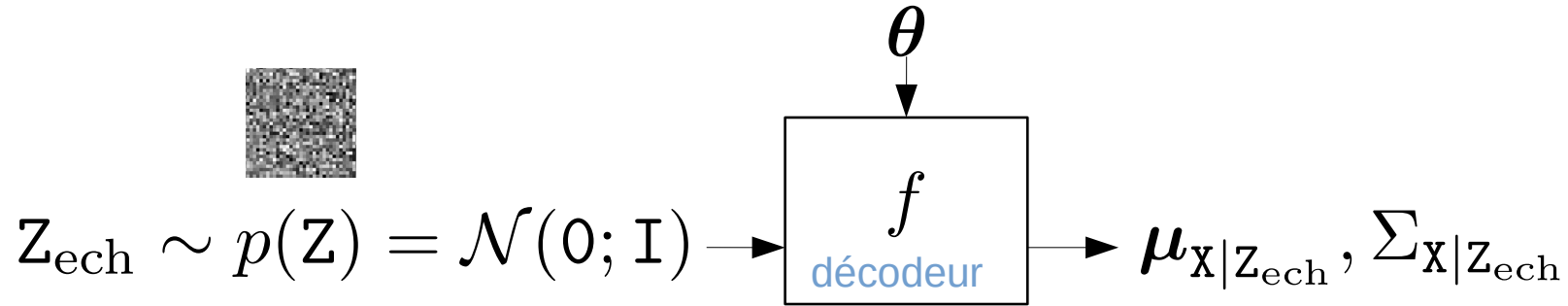
Réseau
inversible



Auto-encodeur
variationnel



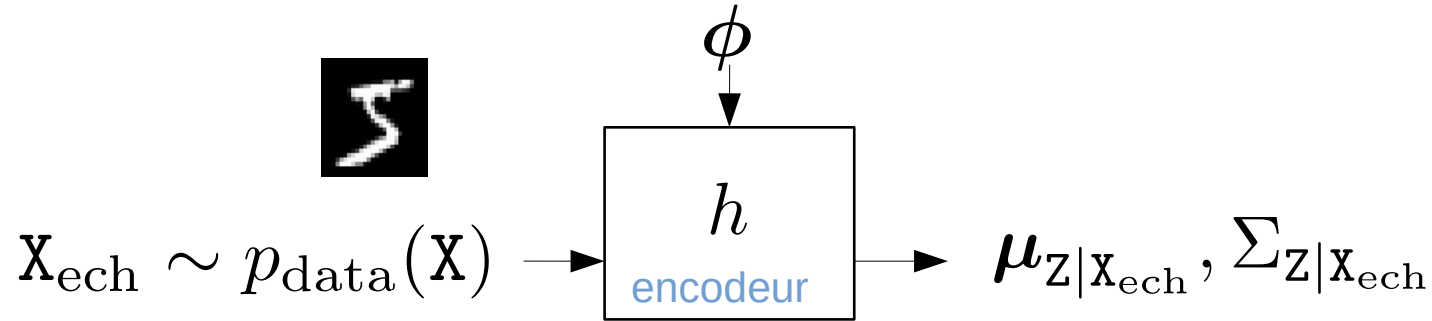
Auto-encodeur variationnel : décodeur



$$p_{\text{m}}(\mathbf{X}|\mathbf{Z}_{\text{ech}}, \boldsymbol{\theta}) = \mathcal{N}(\mu_{\mathbf{X}|\mathbf{Z}_{\text{ech}}}, \Sigma_{\mathbf{X}|\mathbf{Z}_{\text{ech}}}) \xrightarrow{\text{échantillon}} \text{Image}$$

$$p_{\text{m}}(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta}) = p(\mathbf{Z})p_{\text{m}}(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta})$$

Auto-encodeur variationnel : encodeur



$$q_m(Z|X_{\text{ech}}, \phi) = \mathcal{N}(\mu_{Z|X_{\text{ech}}}, \Sigma_{Z|X_{\text{ech}}}) \xrightarrow{\text{échantillon}} \text{Image}$$

$$q_m(X, Z|\phi) = p_{\text{data}}(X)q_m(Z|X, \phi)$$

Auto-encodeur variationnel : apprentissage

Objectif : optimiser θ et ϕ pour que

$$p_m(\mathbf{X}, \mathbf{Z}|\theta) = p(\mathbf{Z})p_m(\mathbf{X}|\mathbf{Z}, \theta) \approx q_m(\mathbf{X}, \mathbf{Z}|\phi) = p_{\text{data}}(\mathbf{X})q_m(\mathbf{Z}|\mathbf{X}, \phi)$$

Pour un VAE on choisit généralement la divergence de Kullback-Leibler suivante :

$$KL(q_m(\mathbf{X}, \mathbf{Z}|\phi) || p_m(\mathbf{X}, \mathbf{Z}|\theta))$$

III)

Auto-encodeur variationnel : apprentissage (suite)

$$KL(q_m(\mathbf{X}, \mathbf{Z}|\phi) || p_m(\mathbf{X}, \mathbf{Z}|\theta))$$

$$= \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}, \Sigma_{\mathbf{z}|\mathbf{x}})} \left(-\ln \mathcal{N}(\boldsymbol{\mu}_{\mathbf{x}|\mathbf{z}}, \Sigma_{\mathbf{x}|\mathbf{z}}) \right) \quad \text{“Erreur de reconstruction” :}$$

Incite la sortie du décodeur à ressembler aux X

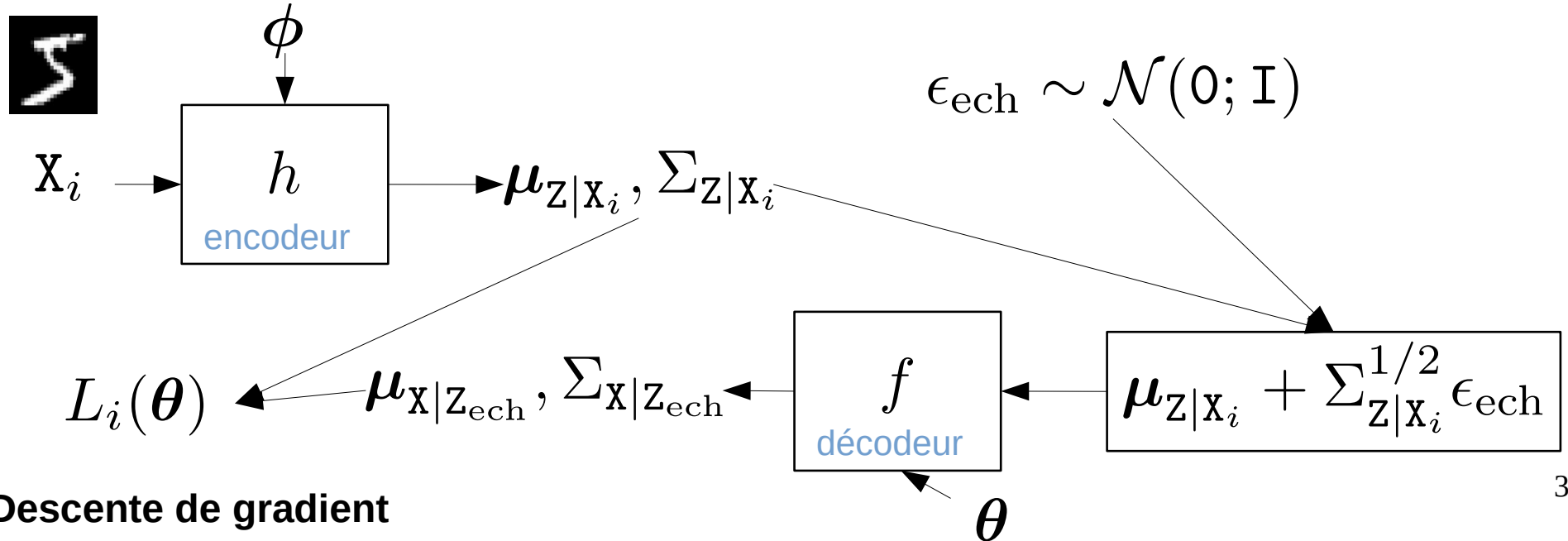
$$+ \mathbb{E}_{p_{\text{data}}(\mathbf{x})} (KL(\mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}, \Sigma_{\mathbf{z}|\mathbf{x}}) || \mathcal{N}(\mathbf{0}, \mathbf{I}))) \quad \text{“Régularisation” :}$$

Incite la sortie de l’encodeur à être centrée réduite → logique car lorsqu’on générera des z, on les tirera selon une gaussienne centrée réduite

$$+ \text{cst}_{\phi, \theta}$$

Auto-encodeur variationnel : apprentissage (suite)

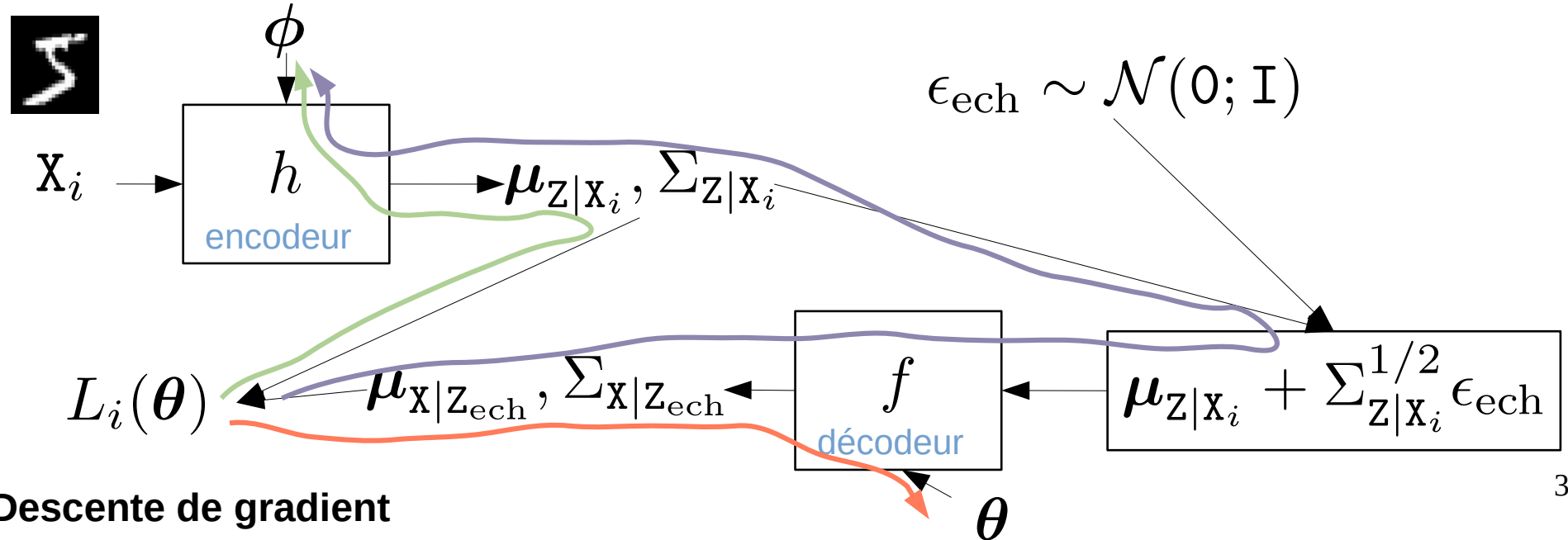
$$\epsilon_{\text{ech}} \sim \mathcal{N}(0; \mathbf{I}) \quad \left| \quad L_i(\theta) = KL(\mathcal{N}(\mu_{\mathbf{Z}|\mathbf{X}_i}, \Sigma_{\mathbf{Z}|\mathbf{X}_i}) || \mathcal{N}(0, \mathbf{I})) \right. \\ \left. - \ln \mathcal{N}(\mu_{\mathbf{X}|\mathbf{Z}}(\mu_{\mathbf{Z}|\mathbf{X}_i} + \Sigma_{\mathbf{Z}|\mathbf{X}_i}^{1/2} \epsilon_{\text{ech}}), \Sigma_{\mathbf{X}|\mathbf{Z}}(\mu_{\mathbf{Z}|\mathbf{X}_i} + \Sigma_{\mathbf{Z}|\mathbf{X}_i}^{1/2} \epsilon_{\text{ech}})) \right.$$



Descente de gradient

Auto-encodeur variationnel : apprentissage (suite)

$$\epsilon_{\text{ech}} \sim \mathcal{N}(0; \mathbf{I}) \quad \left| \quad L_i(\theta) = KL(\mathcal{N}(\mu_{\mathbf{z}|\mathbf{x}_i}, \Sigma_{\mathbf{z}|\mathbf{x}_i}) || \mathcal{N}(0, \mathbf{I})) \right. \\ \left. - \ln \mathcal{N}(\mu_{\mathbf{x}|\mathbf{z}}(\mu_{\mathbf{z}|\mathbf{x}_i} + \Sigma_{\mathbf{z}|\mathbf{x}_i}^{1/2} \epsilon_{\text{ech}}), \Sigma_{\mathbf{x}|\mathbf{z}}(\mu_{\mathbf{z}|\mathbf{x}_i} + \Sigma_{\mathbf{z}|\mathbf{x}_i}^{1/2} \epsilon_{\text{ech}})) \right.$$



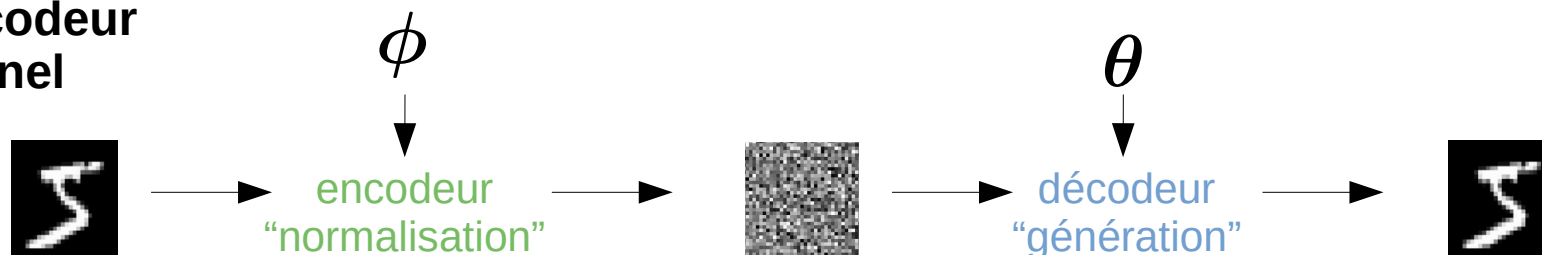
Auto-encodeur variationnel : avantages et inconvénients

- + capable d'échantillonner efficacement
- + pas de contrainte particulière sur l'architecture
- + taille de Z non contrainte par celle de X
- pas d'expression exacte de la (log-)probabilité d'une donnée
- impossible d'apprendre par maximum de vraisemblance (mais borne inférieure quand même)
- contrainte sur la forme des distributions conditionnelles (exemple : covariance diagonale pour avoir un calcul rapide)

IV) Réseau antagoniste (GAN)

Réseau antagoniste : idée générale

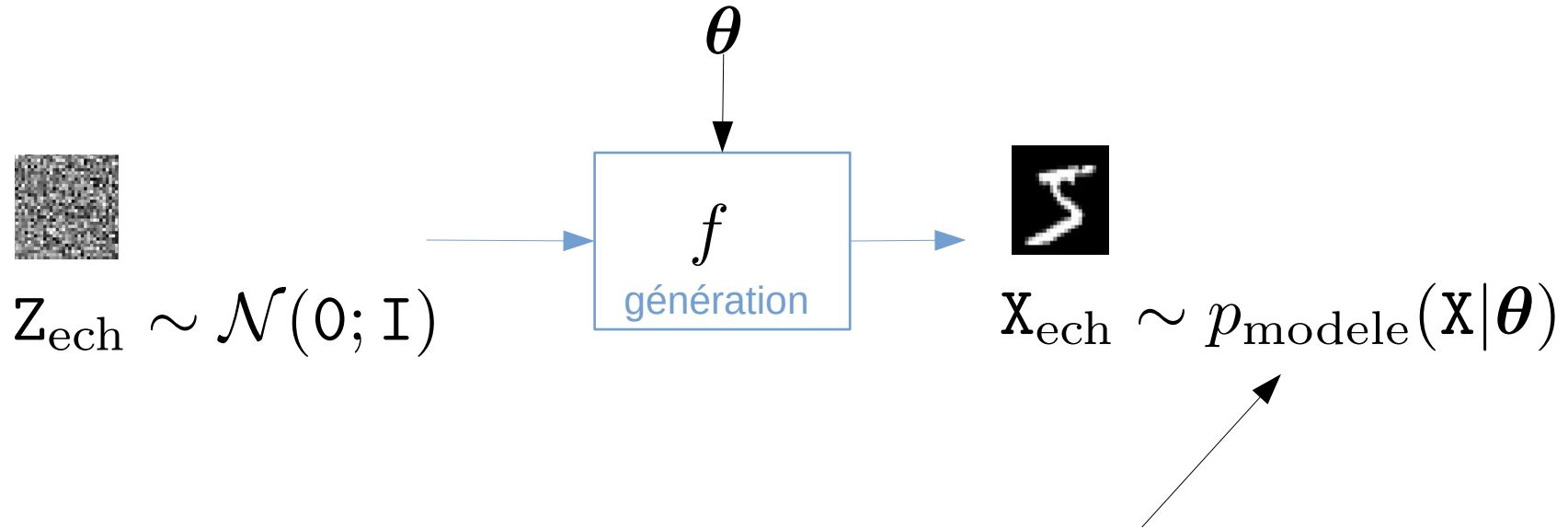
Auto-encodeur
variationnel



Réseau
antagoniste

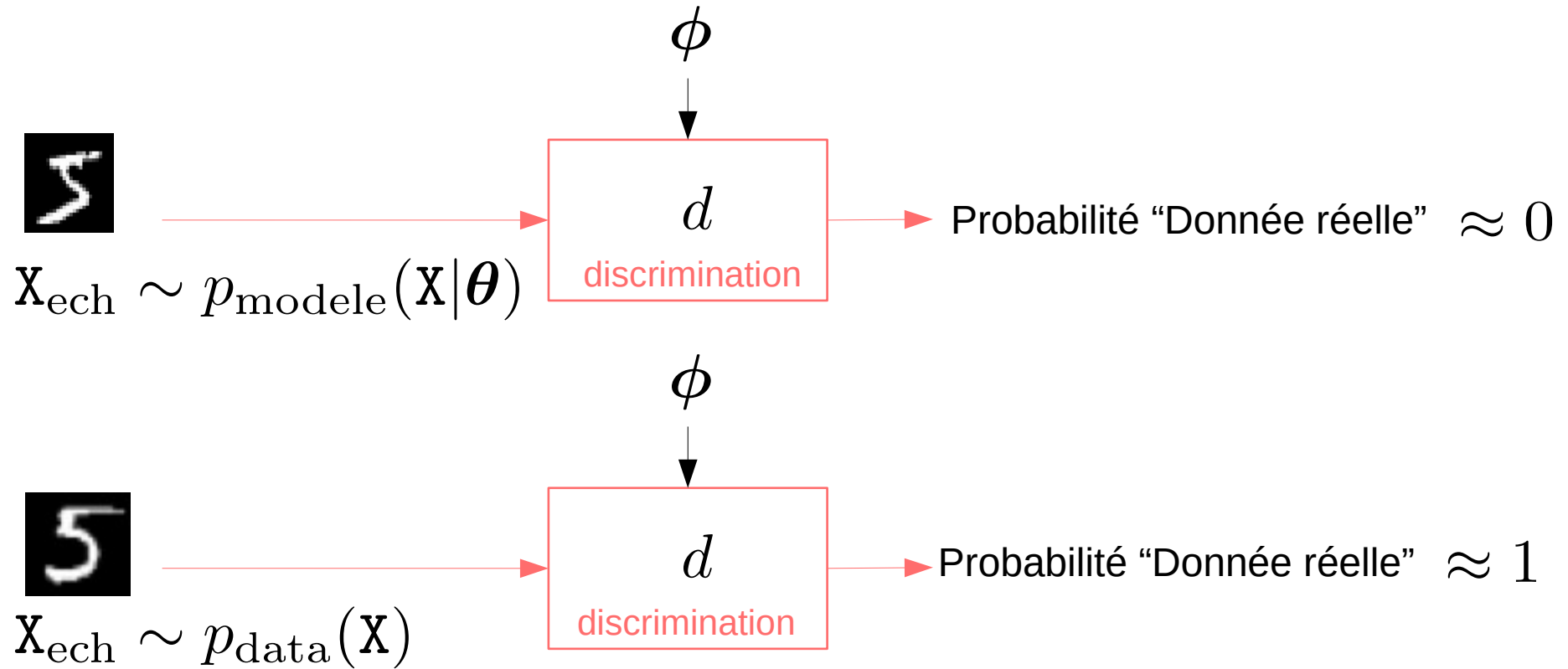


Réseau antagoniste : générateur



Définition implicite, le générateur n'est pas inversible

Réseau antagoniste : discriminateur



Réseau antagoniste : apprentissage

Objectif : optimiser θ et ϕ pour que $p_{\text{modele}}(\mathbf{X}|\theta) \approx p_{\text{data}}(\mathbf{X})$

Pour un GAN on choisit généralement une divergence de la famille suivante :

$$D(p_m(\mathbf{X}|\theta) || p_{\text{data}}(\mathbf{X})) = \max_{\phi} \mathbb{E}_{p_m(\mathbf{x}|\theta)}(h'(\mathbf{X}, \phi)) - \mathbb{E}_{p_{\text{data}}(\mathbf{x})}(h''(\mathbf{X}, \phi))$$

Exemple :

$$\max_{\phi} \mathbb{E}_{p_m(\mathbf{x}|\theta)}(\ln(1 - d(\mathbf{X}, \phi))) + \mathbb{E}_{p_{\text{data}}(\mathbf{x})}(\ln(d(\mathbf{X}, \phi)))$$

Maximum atteint quand $d(\mathbf{X}, \phi) = \frac{p_{\text{data}}(\mathbf{X})}{p_{\text{data}}(\mathbf{X}) + p_m(\mathbf{X}|\theta)}$

Si $p_m = p_{\text{data}}$
Valeur div. $-\ln(4)$

Réseau antagoniste : apprentissage (suite)

$$\begin{aligned} & \min_{\boldsymbol{\theta}} D(p_m(\mathbf{X}|\boldsymbol{\theta})||p_{\text{data}}(\mathbf{X})) \\ &= \min_{\boldsymbol{\theta}} \max_{\phi} \mathbb{E}_{p_m(\mathbf{x}|\boldsymbol{\theta})} (\ln(1 - d(\mathbf{X}, \phi))) + \mathbb{E}_{p_{\text{data}}(\mathbf{x})} (\ln(d(\mathbf{X}, \phi))) \\ &= \min_{\boldsymbol{\theta}} \max_{\phi} \mathbb{E}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} (\ln(1 - d(f(\mathbf{Z}, \boldsymbol{\theta}), \phi))) + \mathbb{E}_{p_{\text{data}}(\mathbf{x})} (\ln(d(\mathbf{X}, \phi))) \end{aligned}$$

Problème minmax → potentiellement beaucoup plus difficile à optimiser qu'un VAE ou qu'un NF

Réseau antagoniste : apprentissage (suite)

En pratique on alterne :

Un pas gradient pour optimiser ϕ

$$L'_i(\phi) = -\ln(1 - d(f(\mathbf{Z}_{\text{ech}}, \boldsymbol{\theta}), \phi)) - \ln(d(\mathbf{X}_i, \phi))$$

Objectif : Améliorer les paramètres du discriminateur afin qu'il prédise une valeur proche de 0 pour une donnée générée et proche de 1 pour une donnée réelle

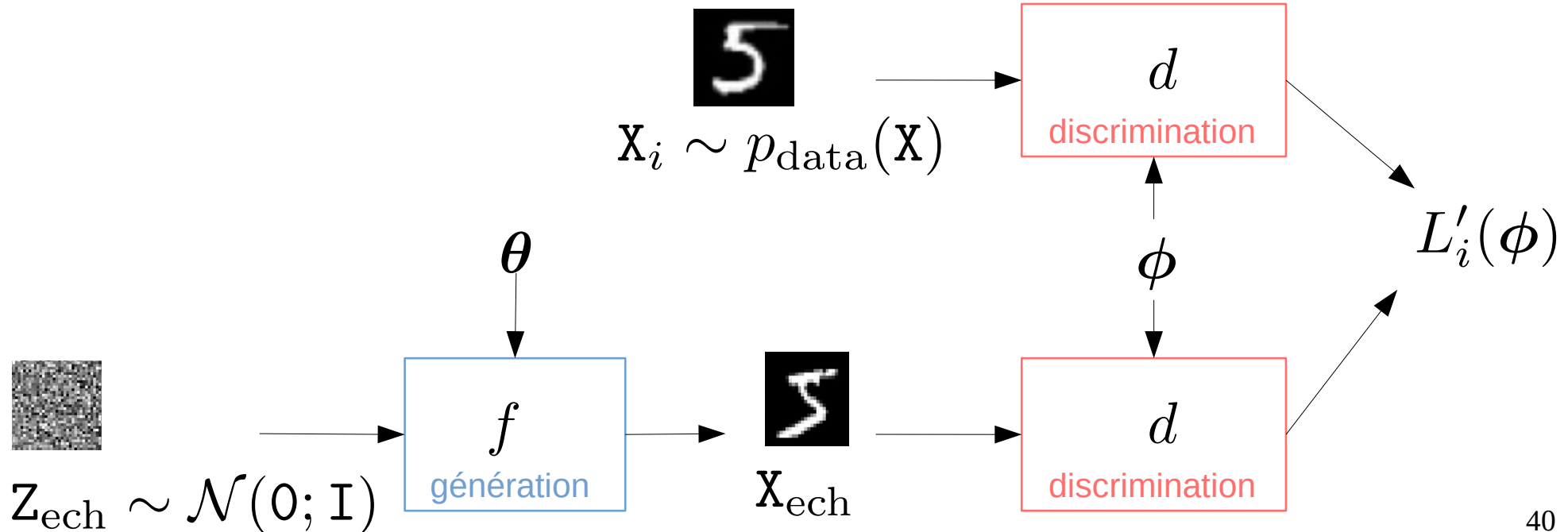
Suivi d'un pas de gradient pour optimiser $\boldsymbol{\theta}$

$$L''_i(\boldsymbol{\theta}) = \ln(1 - d(f(\mathbf{Z}_{\text{ech}}, \boldsymbol{\theta}), \phi)) \longrightarrow \boxed{\text{En pratique remplacé par : } -\ln(d(f(\mathbf{Z}_{\text{ech}}, \boldsymbol{\theta}), \phi))}$$

Objectif : Améliorer les paramètres du générateur afin que le discriminateur se trompe et prédise une valeur proche de 1

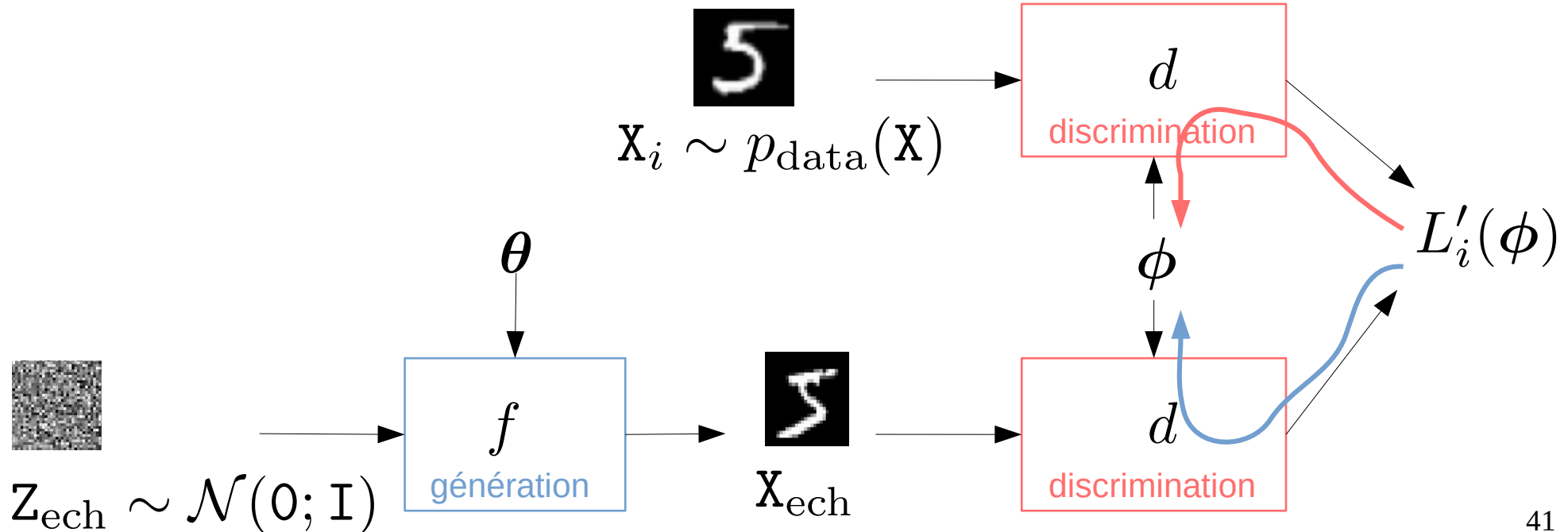
Réseau antagoniste : apprentissage (suite)

Pas de gradient sur les paramètres du discriminateur



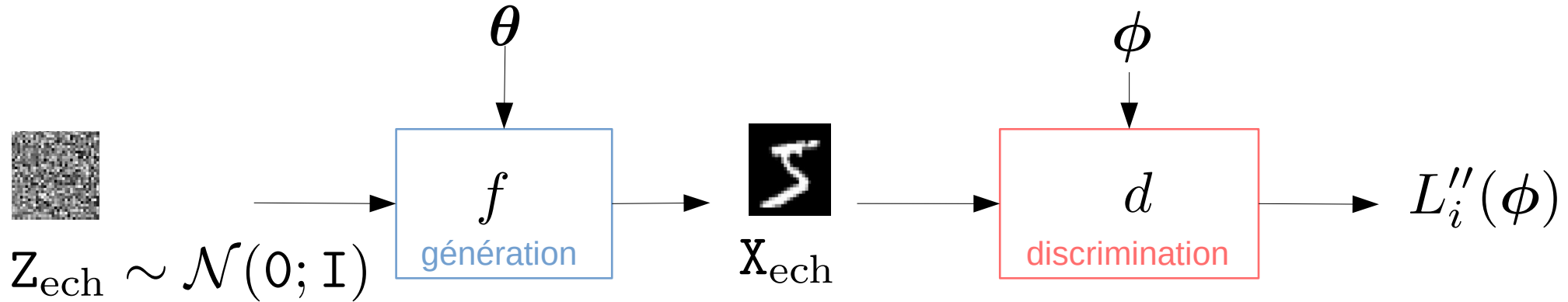
Réseau antagoniste : apprentissage (suite)

Pas de gradient sur les paramètres du discriminateur



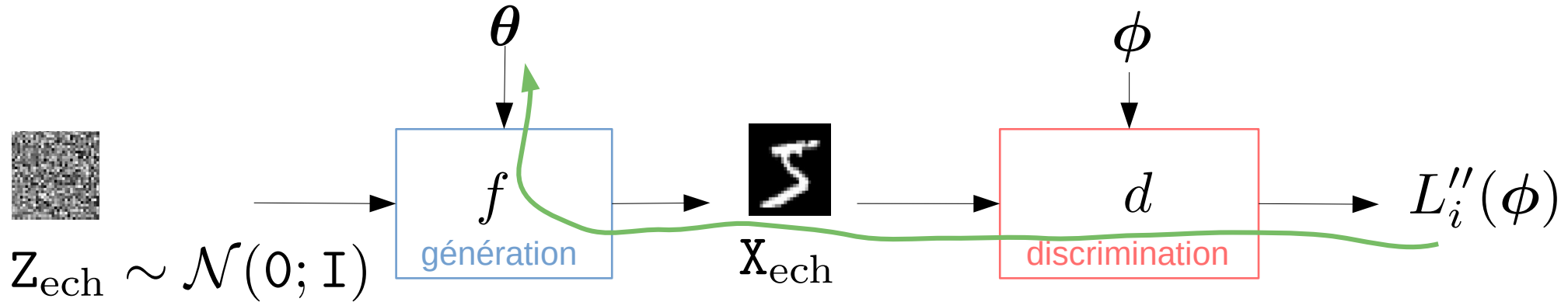
Réseau antagoniste : apprentissage (suite)

Pas de gradient sur les paramètres du générateur



Réseau antagoniste : apprentissage (suite)

Pas de gradient sur les paramètres du générateur



Réseau antagoniste : avantages et inconvénients

- + capable d'échantillonner efficacement
- + pas de contrainte particulière sur l'architecture
- + taille de Z non contrainte par celle de X
- pas d'expression de la (log-)probabilité d'une donnée
- problème minmax difficile à optimiser
 - ** problèmes de convergence
 - ** échantillonnage uniquement d'une partie de la distribution ("mode collapse")

