

# Rapport Projet de Synthèse

---

## **Une utilisation de la méthode dite de *Contrastive Divergence* pour l'apprentissage statistique**

**Axel Flatrès--Berthelot & Baptiste Doyen**

**(Binôme 97 / PAG 2)**

**Juin 2016**



**CentraleSupélec**

# Sommaire

---

## Introduction

### **Partie I** : Approche théorique et fondements mathématiques

- 1) Le modèle énergétique
- 2) L'apprentissage par *Contrastive Divergence*
- 3) L'approximation de *Geoffrey E. Hinton*
- 4) La méthode du *Gibbs Sampling*

### **Partie II** : Implémentation de la solution

- 1) Gaussienne unique
- 2) Mélange de gaussiennes

### **Partie III** : Résultats expérimentaux

## Conclusion et remerciements

# Introduction

---

En termes de méthodologie statistique, il est d'usage de recourir à l'estimateur de maximum de vraisemblance pour estimer les paramètres liés à une distribution aléatoire. En effet, celui-ci permet de disposer d'un estimateur sans biais de manière simple et utile:

- Simple : il suffit de minimiser la vraisemblance  $L(X, \theta)$  et  $\theta_{MV} = \text{argmax}_{\theta} L(X, \theta)$
- Utile: il est asymptotiquement consistant, permet de disposer de l'information de Fisher et est optimal (dans la classe des estimateurs sans biais)

Néanmoins, outre ses qualités théoriques, l'estimateur de maximum de vraisemblance ne présente pas les mêmes avantages en ce qui concerne son calcul pratique.

En effet, si l'on considère une distribution de loi de probabilité de la forme suivante :

$$P(X, \theta) = \frac{1}{Z(\theta)} \exp(X^T \theta X + B^T X)$$

Le calcul de  $\theta_{MV}$  implique celui de  $Z(\theta)$ .

Or, il se peut que  $Z(\theta)$  ne soit pas calculable en pratique (par exemple : s'il s'agit d'une somme contenant un nombre exponentiel de termes).

D'où la nécessité d'une autre méthode pour estimer en pratique le paramètre  $\theta$

L'objet de notre projet de synthèse dont ce texte en est la présentation est d'étudier cette méthode afin de la mettre en place algorithmiquement sur des cas simples (le cas gaussien par exemple).

# Partie I : Théorie

---

## 1) Le modèle énergétique

Afin de réaliser une inférence statistique, c'est-à-dire le passage entre des observations empiriques à la modélisation de ces observations à travers une loi de probabilité, nous allons choisir de toujours utiliser une loi de probabilité de même forme.

Cette forme est la suivante et est inspirée des résultats de physique statistique (statistique de Maxwell-Boltzmann) :

$$\mathbf{P}(\mathbf{X}, \mathbf{W}) = \frac{1}{\mathbf{Z}(\boldsymbol{\theta})} \exp(\mathbf{X}^T \mathbf{W} \mathbf{X})$$

Nous appellerons ce modèle, le modèle énergétique. L'énergie associée a pour expression :

$$\mathbf{E}(\mathbf{X}, \mathbf{W}) = -\mathbf{X}^T \mathbf{W} \mathbf{X}$$

Ainsi, l'énergie associée à  $\mathbf{X}$  est élevée si la probabilité correspondante  $\mathbf{P}(\mathbf{X}, \mathbf{W})$  est faible et vice-versa.

Ce qui est cohérent si l'on interprète ce résultat avec une approche plus « physique » : les configurations les plus probables sont celles qui requièrent le moins d'énergie. Et elles le sont d'autant plus, que l'énergie associée est minimale.

L'un des points forts de ce modèle est son caractère universel : il obéit à un principe selon lequel le plus probable correspond à une minimisation de l'« énergie » nécessaire à sa réalisation.

## 2) L'apprentissage par Contrastive Divergence

**Hypothèse :** on suppose que la distribution aléatoire observée obéit à la règle du modèle énergétique et possède donc la densité associée.

Dans ce cas, on pose

$$a) \mathbf{p}_0 = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n)$$

La distribution expérimentale, c'est-à-dire celle à partir de laquelle l'algorithme « apprend ».

$$b) \mathbf{p}_\infty = \mathbf{p}(\mathbf{W})$$

(La distribution finale, c'est-à-dire celle à partir de laquelle découle l'observation. La théorie qui existerait derrière l'expérimental.)

Au lieu de minimiser la vraisemblance, la méthode dite de Contrastive Divergence se propose de minimiser la différence entre deux divergences, appelées des divergences de Kullback-Leibler :

$$CD_n(\mathbf{W}) = KL(\mathbf{p}_0 \parallel \mathbf{p}_\infty) - KL(\mathbf{p}_n \parallel \mathbf{p}_\infty)$$

Où 
$$KL(\mathbf{p}_0 \parallel \mathbf{p}_\infty) = \sum_x \mathbf{p}_0(x) \log \frac{\mathbf{p}_0(x)}{\mathbf{p}(x, \mathbf{W})}$$

La divergence de Kullback-Keiber entre deux distributions de probabilité mesure la distance entre ces deux distributions (attention : elle est non-symétrique).

Afin de déterminer  $\hat{\mathbf{W}} = \mathit{argmin}_{\mathbf{W}} CD_n(\mathbf{W})$ , on utilise une méthode dite de descente de gradient.

On définit la suite  $\mathbf{W}_n$  de la sorte :

$$\mathbf{W}_{n+1} = \mathbf{W}_n + \frac{\partial CD_n(\mathbf{W})}{\partial \mathbf{W}} \bigg|_{\mathbf{W}_n}$$

Cette suite converge vers la valeur souhaitée, à savoir  $\hat{\mathbf{W}}$ .

Reste à déterminer  $\frac{\partial CD_n(\mathbf{X}, \mathbf{W})}{\partial \mathbf{W}} \bigg|_{\mathbf{W}_n}$  : pour cela nous allons utiliser l'approximation de Hinton.

### 3) L'approximation de *Geoffrey E. Hinton*

D'après ce qui précède :

(Par souci de simplification on prend  $\frac{\partial}{\partial W} = \frac{\partial}{\partial W_{W_n}}$  )

$$\frac{\partial CD_n(X, W)}{\partial W} = \frac{\partial KL(p_0 \parallel p_\infty)}{\partial W} - \frac{\partial KL(p_n \parallel p_\infty)}{\partial W}$$

Or  $\frac{\partial KL(p_0 \parallel p_\infty)}{\partial W} = \langle \frac{\partial \log p_0}{\partial W} \rangle_{p_0} - \langle \frac{\partial \log p_\infty}{\partial W} \rangle_{p_0}$  ( $p_0$  indépendant de  $W_n$ )

De même :  $\frac{\partial KL(p_n \parallel p_\infty)}{\partial W} = \langle \frac{\partial \log p_n}{\partial W} \rangle_{p_n} - \langle \frac{\partial \log p_\infty}{\partial W} \rangle_{p_n}$

Et :  $-\langle \frac{\partial \log p_\infty}{\partial W} \rangle_{p_0} + \langle \frac{\partial \log p_\infty}{\partial W} \rangle_{p_n}$  est considéré comme négligeable. (Hinton)

D'où :  $\frac{\partial CD_n(X, W)}{\partial W} \approx \langle \frac{\partial \log p_0}{\partial W} \rangle_{p_0} - \langle \frac{\partial \log p_n}{\partial W} \rangle_{p_n} \approx \langle \mathbf{X}^T \mathbf{X} \rangle_{p_0} - \langle \mathbf{X}^T \mathbf{X} \rangle_{p_n}$

$p_n$  n'est pas indépendant de  $W_n$  mais par effet de différence le calcul de la fonction de partition est ainsi évité par différence des deux quantités ci-dessus.

Afin de déterminer  $\langle \mathbf{X}^T \mathbf{X} \rangle_{p_n}$ , on utilise l'échantillonnage de Gibbs.

## 4) L'échantillonnage de Gibbs

Une forme **sans variables cachées**, à partir des lois jointes :  
on la nomme Echantillonnage de Gibbs cyclique et elle repose sur la relation suivante.

Soit  $(\mathbf{X}^j)$  une famille de vecteurs de dimension  $n$  telle que  $\mathbf{X}^j = (x_1^j, \dots, x_n^j)$ .

$$x_i^j \sim P(x_i^j | \mathbf{x}_{-i}^j) = \frac{1}{1 + e^{-\sum_{j \neq i} \theta_{i,j} x_j^j}}$$

avec  $\mathbf{x}_{-i}^j$  le vecteur  $\mathbf{X}^j$  ôté de sa composante  $x_i^j$ .

Etant donné un vecteur, suite à l'échantillonnage de ses composantes grâce à la loi jointe ci-dessus, on peut réaliser une descente de gradient pour obtenir une nouvelle matrice des paramètres.

---

# Partie II : implémentation de la solution

---

## 1) Gaussienne unique

**Objectif** : déterminer les paramètres d'une gaussienne à savoir sa moyenne  $\mu$  et sa variance  $\sigma^2$ .

**Entrées** :  $m$  valeurs générées aléatoirement à partir de  $\mathcal{N}(\mu, \sigma^2)$ .

**Algorithme** :

```
%taux d'apprentissage
epsilon = 10^(-3);

n = 10;
m = 1000;

%génération des training examples

mu = 45;
sigma2 = 10;

Y = normrnd(mu, sigma2, [1 m]);

%conversion en binaire des données
Y_int = floor(Y);
Y_binary = de2bi(Y_int, n);

X1 = zeros(n, m);
X2 = zeros(n, m);

%Initialisation des paramètres
W0 = zeros(n, n);
```



```
%Début entraînement du modèle sur les training examples
```

```
for i =1:m
```

```
    X1(:,i) = Y_binary(i,:);
```

```
    for d = 1:n
```

```
        s1 = 0;
```

```
        for c = 1:n
```

```
            if c ~= d
```

```
                s1 = s1 + W0(d,c)*X1(c,i);
```

```
            end;
```

```
        end;
```

```
        p1 = (1/(1+exp(-s1)));
```

```
        p0 = 1-p1;
```

```
        X2(d,i) = randsample(0:1,1,true,[p0 p1]);
```

```
    end;
```

```
%Début descente de Gradient
```

```
for a=1:n
```

```
    for b=1:n
```

```
        W0(a,b) = W0(a,b) +  
        epsilon*0.5*(X1(a,i)*X1(b,i) -  
        X2(a,i)*X2(b,i));
```

```
    end;
```

```
end;
```

```
%Fin descente de Gradient
```

```
end;
```

```
%Fin entraînement du modèle sur les training examples
```

```
%Début de la comparaison entre les résultats obtenus et le modèle  
théorique générateur.
```

```
k = 100;  
X = zeros(1,k);  
result = zeros(1,k);
```

```
for i = 1:k  
    X(i)=i;  
end;
```

```
X_binary = de2bi(X,n);
```

```
for i = 1:k  
    result(i) =  
        exp(0.5*X_binary(i,:)*W0*transpose(X_binary(i,:)));  
end;
```

```
%Détermination de la constante de normalisation  
Q = trapz(X,result);
```

```
%Densité théorique  
for i = 1:k  
    Xnormal(i) = (1/sqrt(sigma2*2*pi))*exp(-0.5*(X(i)-mu)*(X(i)-  
mu)/sigma2);  
end;
```

```
%Affichage des résultats  
figure  
plot(X,result/Q, '*');  
hold on  
plot(X,Xnormal, '*');
```

## 2) Deuxième cas : reconnaître un mélange de gaussienne (*Gaussian Mixture*)

**Objectif** : déterminer les paramètres d'un mélange de gaussiennes à savoir les moyennes  $\mu_1, \mu_2, \dots, \mu_n$ , les variances  $\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2$  et les poids associées à chaque gaussienne  $p_1, p_2, \dots, p_n$

**Entrées** : m valeurs générées aléatoirement à partir de  $\sum_i p_i \mathcal{N}(\mu_i, \sigma_i^2)$ .

### Algorithme :

L'algorithme est le même que celui qui précède à la différence de :

```
%génération des training examples

d = 1;
mu = [20;50];

k = [5 5];

sigma = zeros(1,d,length(k));

for i = 1:length(k)
    sigma(:, :, i) = k(i);
end;

%poids des gaussiennes
p = [1 2];

%génération de l'objet gmdistribution
obj = gmdistribution(mu,sigma,p);

Y = random(obj,m);

%Densité théorique

Xnormal = pdf(obj,transpose(X)) ;
```

# Partie III : Résultats

---

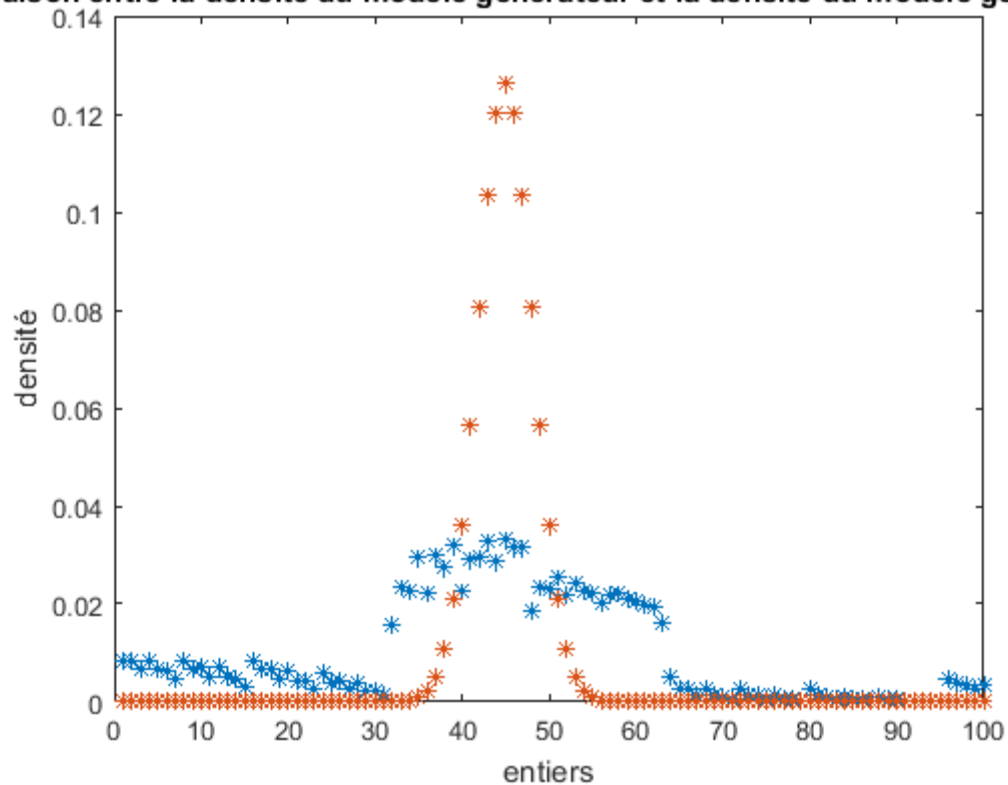
## 1) Gaussienne unique

Comparaison entre la densité du modèle générateur des valeurs initiales et la densité du modèle généré par descente du gradient.

Paramètres :  $(\mu, \sigma^2) = (45, 5)$

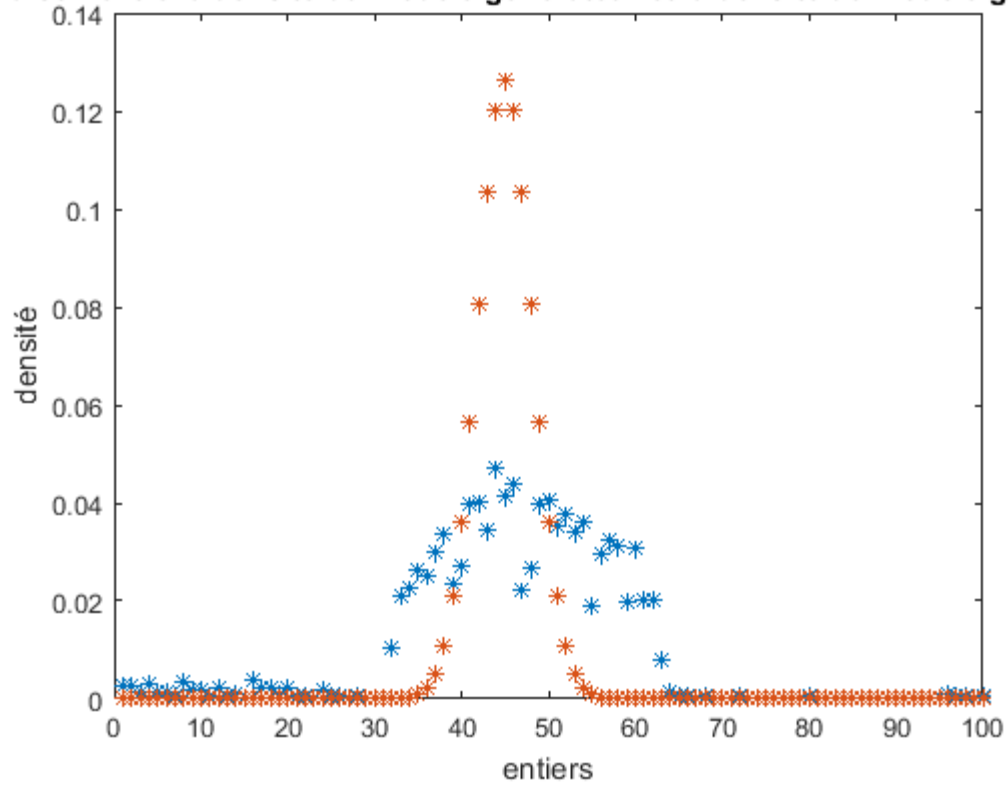
**m = 10 000**

Comparaison entre la densité du modèle générateur et la densité du modèle généré.



**m = 50 000**

comparaison entre la densité du modèle générateur et la densité du modèle généré



**Interprétation :**

- Les valeurs sont symétriques par rapport à la valeur moyenne
- On distingue un pic central autour de la valeur moyenne.
- Plus m augmente, plus la figure est satisfaisante

## 2) Mélange de gaussiennes

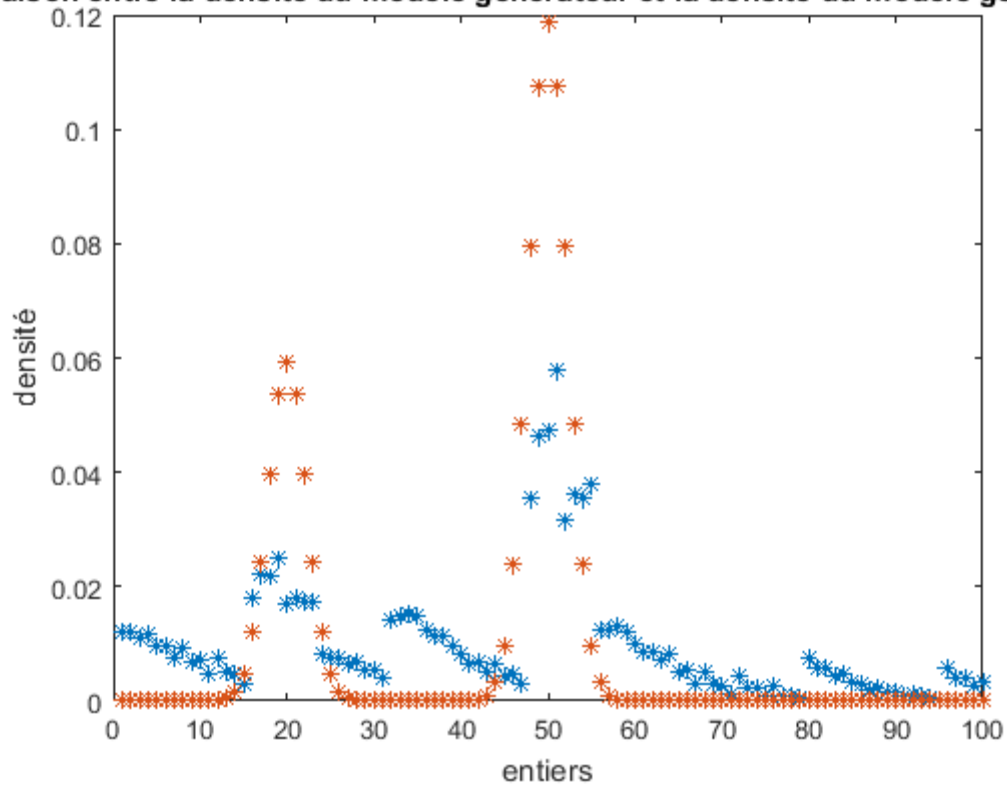
Comparaison entre la densité du modèle générateur des valeurs initiales et la densité du modèle généré par descente du gradient.

Paramètres :  $(\mu_1, \sigma_1^2) = (20, 5)$  et  $(\mu_2, \sigma_2^2) = (50, 5)$

Poids = (1,2)

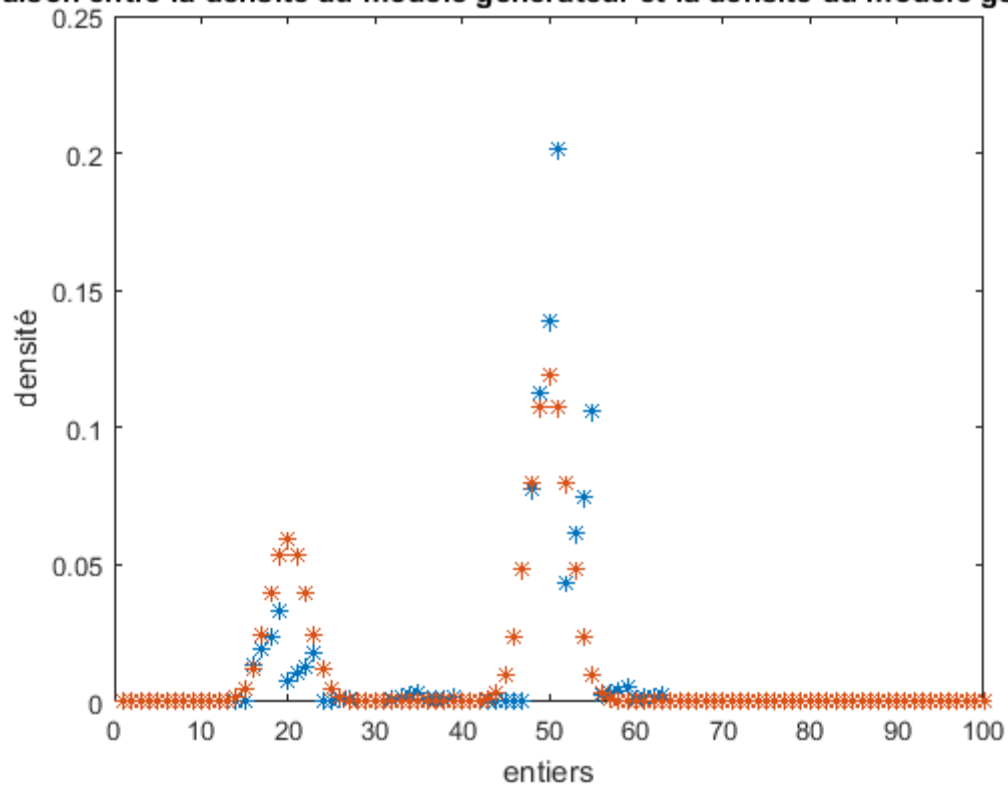
**m = 1000**

Comparaison entre la densité du modèle générateur et la densité du modèle généré.



**m = 5000**

comparaison entre la densité du modèle générateur et la densité du modèle généré.



**Interprétation :**

- Les valeurs sont symétriques par rapport à la valeur moyenne
- On distingue un pic central autour de la valeur moyenne.
- Contrairement au cas précédent, les valeurs sont moins recentrées autour de la valeur moyenne (surtout pour le deuxième pic)
- Plus m augmente, plus la figure est satisfaisante

# Conclusion

---

Dans le cadre de ce projet, nous nous sommes familiarisés avec Matlab. Cela a été très enrichissant.

De plus, les résultats obtenus semblent plutôt satisfaisants au regard des comparaisons graphiques présentées ci-dessus.

Les valeurs prises ont de plus été réduites à des entiers, ce qui a rajouté de l'imprécision aux résultats obtenus.

Pour conclure, cette méthode peut être qualifiée d'efficace car avec peu de calculs et de lignes de code, on peut retrouver des valeurs satisfaisantes du modèle théorique.

Nous tenons également à remercier chaleureusement notre tuteur M. Pablo Piantanida pour son tutorat et son aide au cours de ce projet.