



# La théorie des probabilités : la logique de la découverte scientifique



33ème Festival d'Astronomie de Fleurance,  
Cours Fil Noir

Florent Leclercq

[www.florent-leclercq.eu](http://www.florent-leclercq.eu)

Institut d'Astrophysique de Paris  
CNRS & Sorbonne Université

7 août 2023

# Le dépôt Github

- [https://github.com/florent-leclercq/Bayes\\_InfoTheory](https://github.com/florent-leclercq/Bayes_InfoTheory)

GitHub, Inc. [US] github.com/florent-leclercq/Bayes\_InfoTheory

Search or jump to... Pull requests Issues Marketplace Explore

florent-leclercq / Bayes\_InfoTheory Unwatch 2 Star 0 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights Settings

Lectures on Bayesian statistics and information theory Edit

Manage topics

39 commits 1 branch 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find File Clone or download

florent-leclercq updated notebooks, corrected error in ABC discrepancy

data	added machine learning
.gitignore	updated gitignore
ABC_discrepancy_effective_likelihood.ip...	updated notebooks, corrected error in ABC discrep
ABC_rejection.ipynb	updated ABC notebooks
ABC_synthetic_likelihood.ipynb	updated notebooks, corrected error in ABC discrepancy

Clone with SSH Use HTTPS

Use an SSH Search Copy Send to My Flow

git@github.com:florent-leclercq/Bayes

Download ZIP

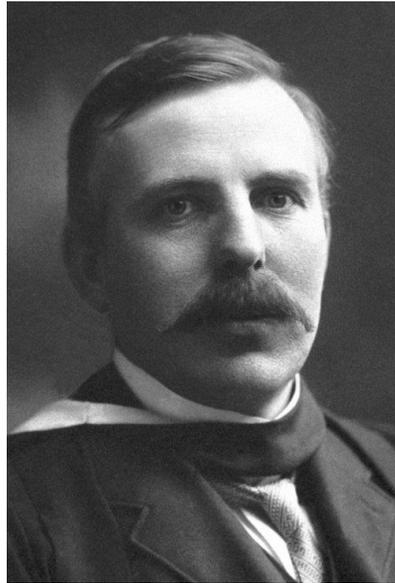
`git clone https://github.com/florent-leclercq/Bayes_InfoTheory.git` (or with SSH)

- Page du cours : <http://florent-leclercq.eu/teaching.php>
- Slides : [http://florent-leclercq.eu/documents/teaching/Cours\\_Fil\\_Noir\\_Fleurance\\_2023.pdf](http://florent-leclercq.eu/documents/teaching/Cours_Fil_Noir_Fleurance_2023.pdf)

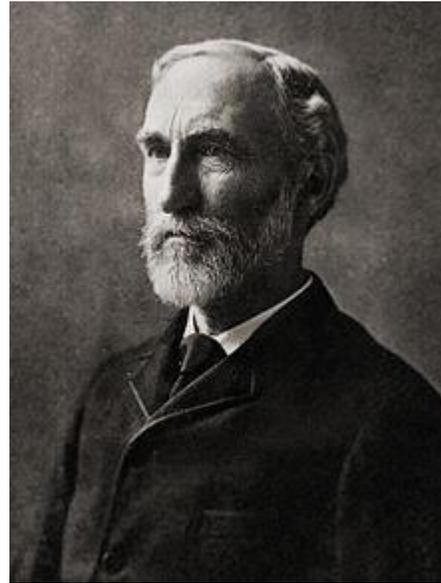
## Introduction : pourquoi faire des statistiques correctement est important

### Un exemple historique : le paradoxe de Gibbs

- « Si votre expérience nécessite des statistiques, vous auriez dû réaliser une meilleure expérience. »



Ernest Rutherford  
(1871-1937)



J. Willard Gibbs  
(1839-1903)

- L'ensemble canonique et l'ensemble grand-canonique de Gibbs, dérivés du principe de l'entropie maximale, **échouent à prédire correctement** les propriétés thermodynamiques des systèmes physiques réels.
- Les entropies prédites sont toujours plus grandes que celles observées... il doit donc exister des **contraintes microphysiques supplémentaires** :
  - Discrétisation des niveaux d'énergie : rayonnement : Planck (1900), solides : Einstein (1907), Debye (1912), Ising (1925), atomes individuels : Bohr (1913)...
  - ... Mécanique quantique : Heisenberg, Schrödinger (1927).

Les premiers indices indiquant le besoin d'une physique *quantique* ont été découverts grâce à des applications apparemment « infructueuses » des statistiques.

# Comment raisonner rationnellement en présence d'incertitudes : La formule de Bayes

Ce n'est (probablement !) pas la bonne personne.

- Comment mesurer une grandeur ?  
Comment vérifier une théorie ? Plus généralement, comment la connaissance progresse-t-elle ?
- La formule de Bayes (1763): une identité mathématique sur la façon dont nous **analysons des observations** et **changeons d'avis** lorsque nous obtenons de nouvelles informations.
- Mais pourquoi l'utiliser ?
  - La formule de Bayes est triviale et dépassée.
  - Elle mesure la **croissance**. Elle dit que nous pouvons apprendre même à partir de données manquantes ou incomplètes, à partir d'approximations, à partir de l'ignorance. Elle va à l'encontre de la conviction que la science exige objectivité et précision.
  - Après la mort de Laplace, elle a été déclarée morte et enterrée.



Thomas Bayes  
(1701-1761)



Richard Price  
(1723-1791)



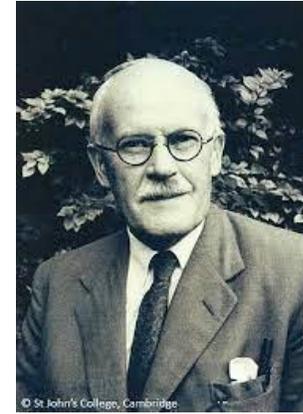
Pierre-Simon de Laplace  
(1749-1827)



Photos prises à Bunhill Fields Burial Ground, City of London, en 2021

## Controverse : fréquentisme versus bayésianisme

- Deux conceptions de la nature des probabilités et des questions scientifiques :
  - Probabilités « objectives » liées à la fréquence de phénomènes aléatoires répétitifs. Questions liées à des expériences déterminées et reproductibles.
  - Probabilités « subjectives » liées à la certitude accordée à une mesure ou à une théorie. Questions liées à des phénomènes et des choix n'impliquant pas l'idée d'une répétition.



Harold Jeffreys  
(1891-18989)



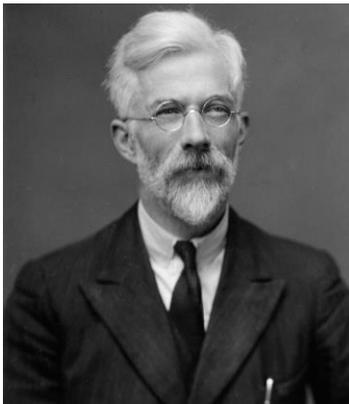
Leonard J. Savage  
(1917-1971)

Fisher publiait parfois des insultes que seul un saint pourrait entièrement pardonner.

Savage 1976, *Annals of Statistics*



Karl Pearson  
(1857-1936)



Ronald Aylmer Fisher  
(1890-1962)



Jerzy Neyman  
(1894-1981)

- Techniques fréquentistes et bayésiennes donnent les mêmes résultats dès qu'on travaille sur des grands échantillons. Ce n'est que sur des petits nombres et des faibles occurrences que l'estimation fréquentiste et l'induction bayésienne diffèrent.

## La théorie qui ne voulait pas mourir

- Et pourtant, la formule de Bayes a aidé à résoudre beaucoup de problèmes pratiques :
  - Prouver l'innocence d'Alfred Dreyfus (Henri Poincaré, 1899-1906),
  - Sauver le système téléphonique Bell de la panique financière (Edward C. Molina, 1907),
  - Prédire les tremblements de terre et les tsunamis (Harold Jeffreys, 1930-1940),
  - Casser le code Enigma de la marine allemande (Alan Turing, 1940-1944),
  - Prouver que la cigarette cause le cancer du poumon (Jerome Cornfield, 1951)...
  - Rechercher un bombe H puis un sous-marin perdus en mer (John P. Craven, 1966-1968)

the theory that would not die  
how bayes' rule cracked the enigma code, hunted down russian submarines & emerged triumphant from two centuries of controversy  
sharon bertsch mcgrayne

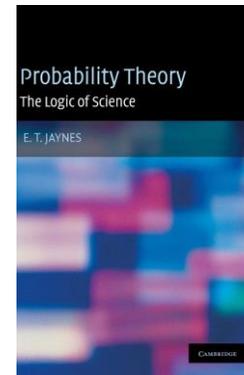
"If you're not thinking like a Bayesian, perhaps you should be."  
—John Allen Paulos, *New York Times Book Review*

Sharon Bertsch McGrayne (2012)

- La bataille scientifique a duré 150 ans, jusqu'à l'arrivée des ordinateurs.

La supériorité des méthodes bayésiennes est aujourd'hui un fait largement démontré dans une centaine de domaines différents. On peut argumenter avec une philosophie ; il n'est pas si facile d'argumenter avec une sortie d'ordinateur, qui nous dit : « Indépendamment de toute votre philosophie, voici les faits de performance réelle ».

Jaynes 2002, *Probability Theory – The logic of science*



Jaynes (2002)



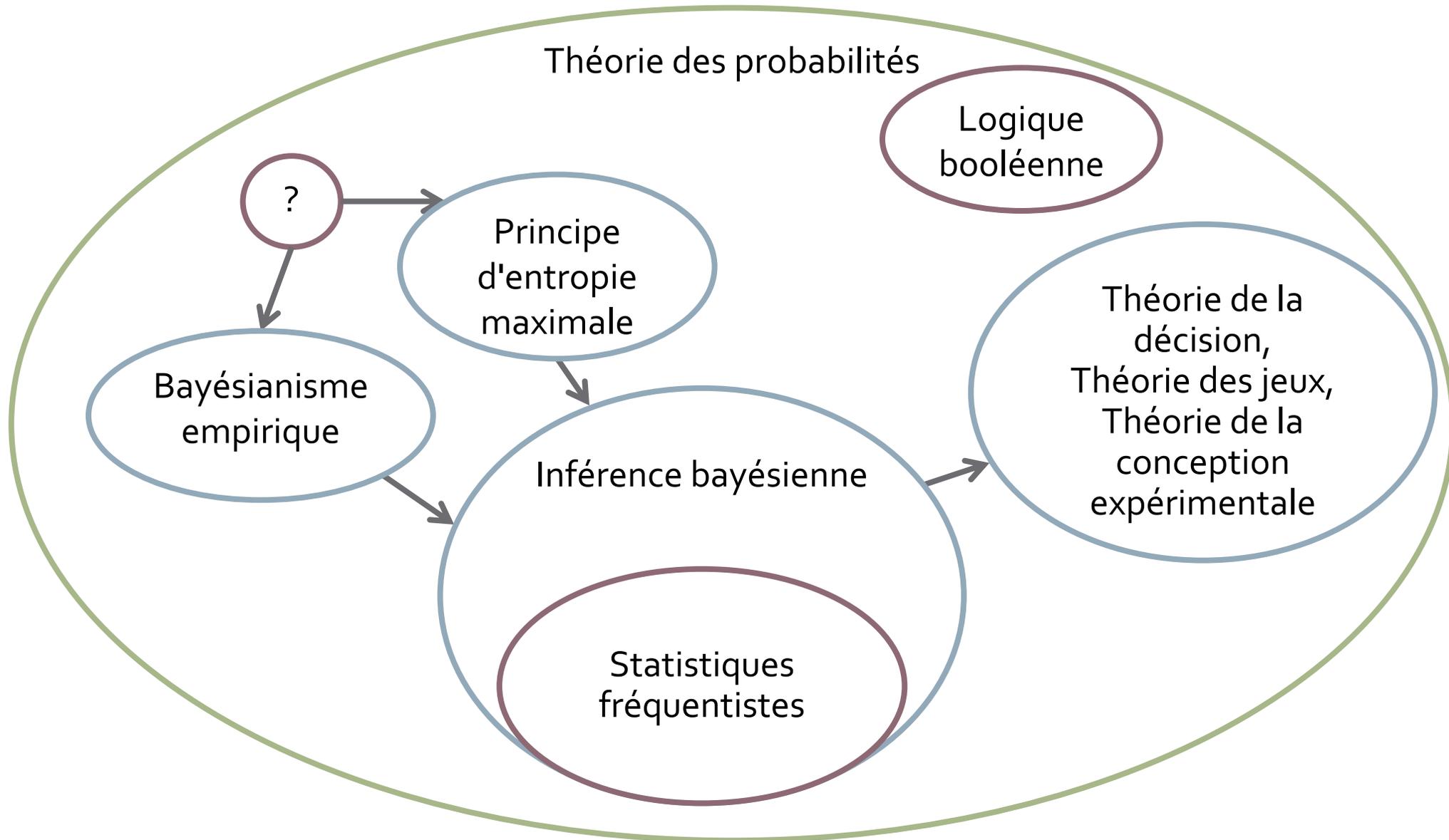
Richard Threlkeld Cox  
(1898-1991)



Edwin Thompson  
Jaynes (1922-1998)

- Théorème de Cox-Jaynes (1946) : la théorie (bayésienne) des probabilités est la seule « logique de la science » possible.

# La « théorie des probabilités » d'après Jaynes : une extension de la logique ordinaire



## Plan – vous pouvez voter pour les sujets que vous souhaitez !

- Principes de base de la théorie des probabilités
- Champs gaussiens aléatoires
- Traitement bayésien du signal et filtrage de Wiener
  - Débruitage bayésien
  - Séparation bayésienne des composantes
- Méthodes Monte Carlo et chaînes de Markov : l'algorithme de Metropolis-Hastings
- Exemple d'inférence de paramètres (cosmologie avec les données des supernovæ)
  - Chaînes de Markov hamiltoniennes
  - Émulateurs
- Exemple de comparaison de modèles (l'éclipse de 1919)
- Priors d'ignorance et le principe de l'entropie maximale
- Théorie de l'information et apprentissage automatique supervisé
- Théorie bayésienne de la décision et théorie bayésienne de la conception expérimentale
- Réseaux bayésiens, modèles bayésiens hiérarchiques et Bayésianisme empirique

## Les principes de base de la théorie des probabilités

- Les probabilités jointes et les probabilités conditionnelles :  $p(A, B) = p(A|B)p(B)$
- Les probabilités marginales :  $p(A) = \int p(A, B) dB$
- La règle du produit (« et » logique) :  $p(AB|C) = p(A|BC) p(B|C)$
- La règle de la somme (« ou » logique) :  $p(A + B|C) = p(A|C) + p(B|C) - p(AB|C)$
- La formule de Bayes (inférence des paramètres d'un modèle) :

$$p(s|d) = \frac{p(d|s) p(s)}{p(d)}$$

Diagram illustrating the components of Bayes' formula:

- posterior:  $p(s|d)$
- vraisemblance:  $p(d|s)$
- prior:  $p(s)$
- évidence:  $p(d)$

- Le facteur de Bayes (comparaison bayésienne de modèles) :  $\mathcal{B}_{12} = \frac{p(d|\mathcal{M}_1)}{p(d|\mathcal{M}_2)}$   
où  $p(d|\mathcal{M}_i) = \int p(d|s_i, \mathcal{M}_i)p(s_i|\mathcal{M}_i) ds_i$

Le ratio des évidences (et non des vraisemblances) prend en compte l'effet « rasoir d'Occam ».

# Champs gaussiens aléatoires

Notebook : <https://github.com/florent->

[leclercq/Bayes\\_InfoTheory/blob/master/GRF\\_and\\_fNL.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/GRF_and_fNL.ipynb)

- Définition : tout vecteur aléatoire  $x$  de distribution de probabilité (pdf) :

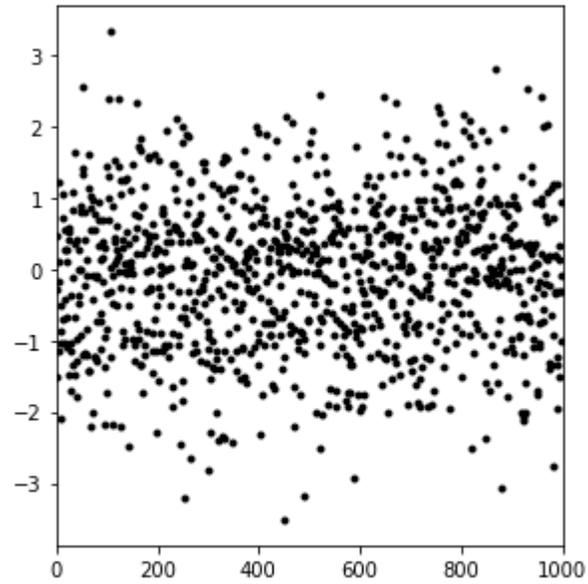
$$p(x|\mu, C) = \frac{1}{\sqrt{|2\pi C|}} \exp \left[ -\frac{1}{2} (x - \mu)^\top C^{-1} (x - \mu) \right]$$

$$-2 \ln p(x|\mu, C) = (x - \mu)^\top C^{-1} (x - \mu) + \ln |2\pi C| \quad \dots \text{c'est tout !}$$

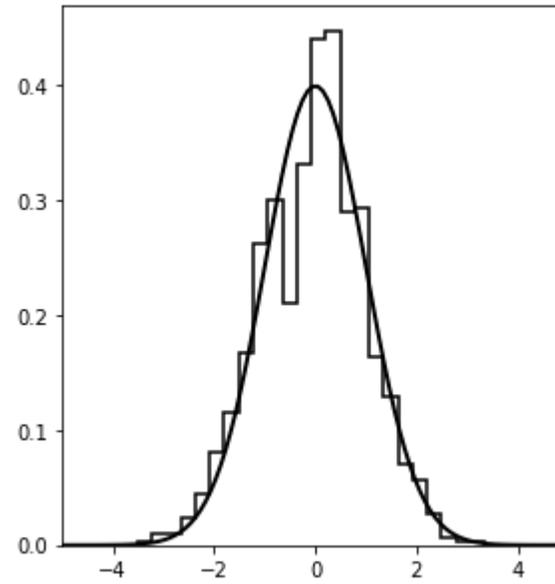
- Génération de champs gaussiens aléatoires :
  - Tirer un vecteur de bruit blanc  $\xi$  (variables gaussiennes non corrélées et de moyenne nulle et variance unitaire)
  - Trouver la racine carrée de la matrice  $C$  :  $\sqrt{C}$  (n'importe quelle matrice satisfaisant cette condition convient)
  - Calculer :  $x = \sqrt{C}\xi + \mu$

# Champs gaussiens aléatoires : exemples

bruit blanc



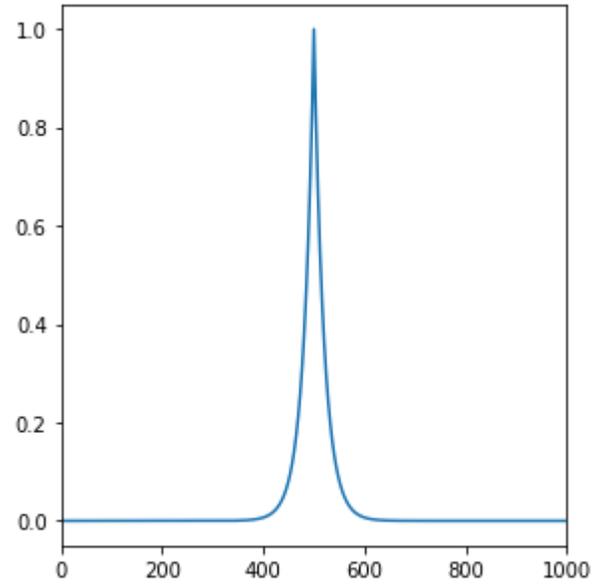
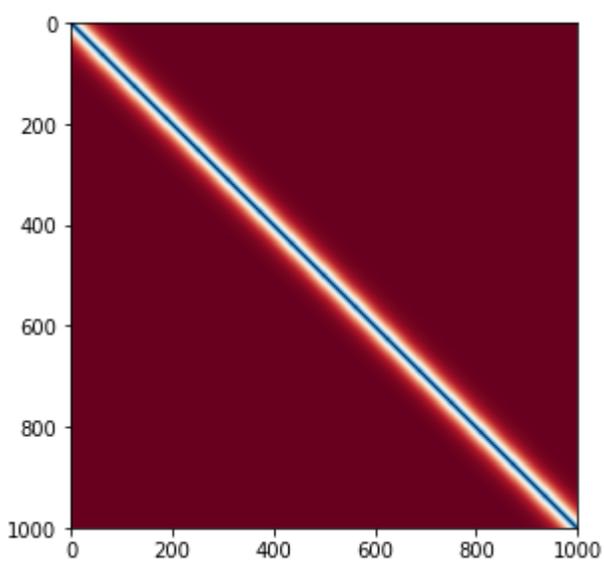
histogramme



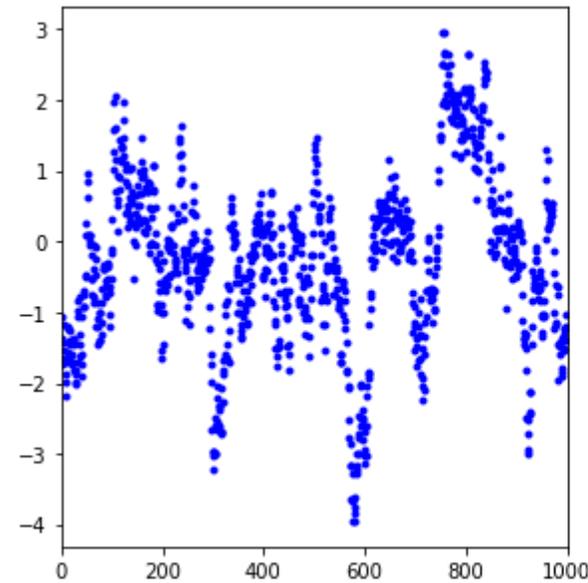
# Champs gaussiens aléatoires : exemples

$$C_{ij} = \exp\left(-\frac{|i-j|}{20}\right)$$

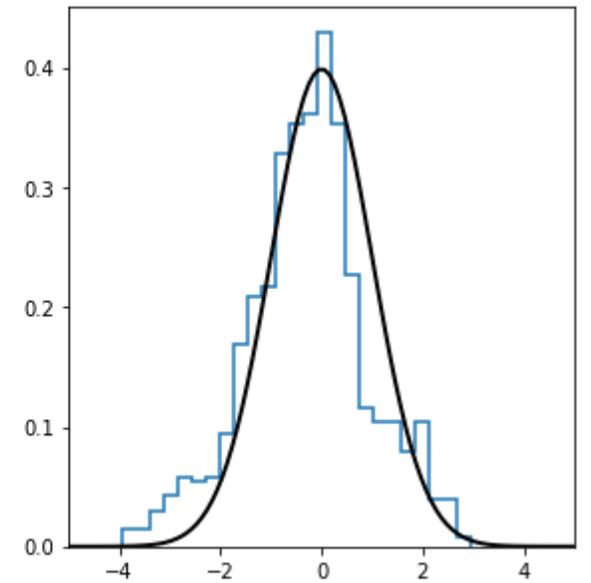
matrice de covariance



champ gaussien



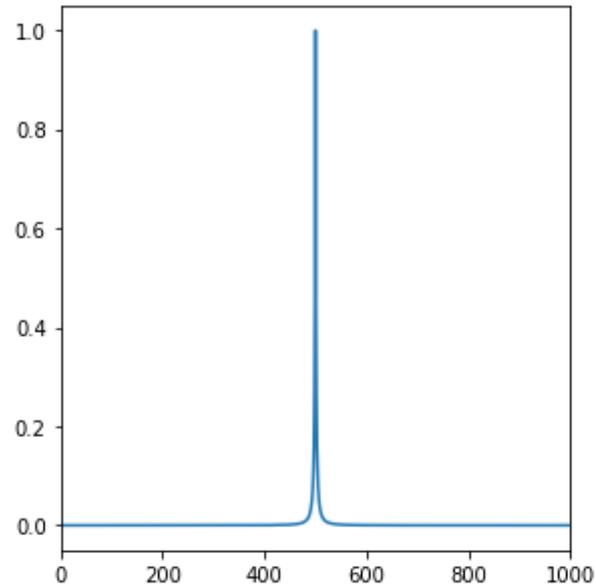
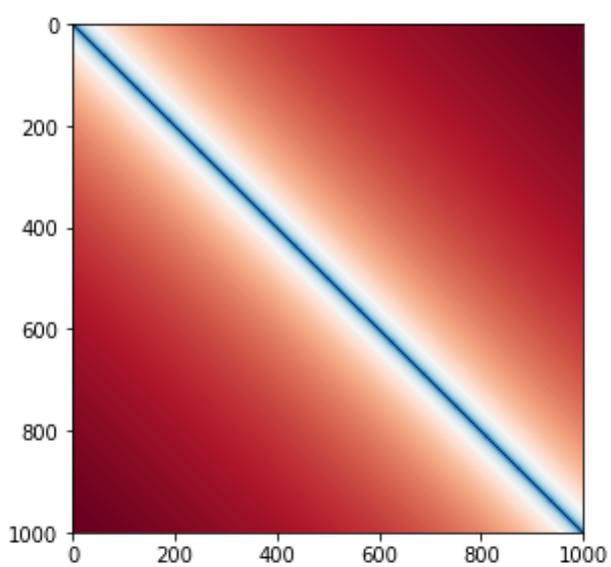
histogramme



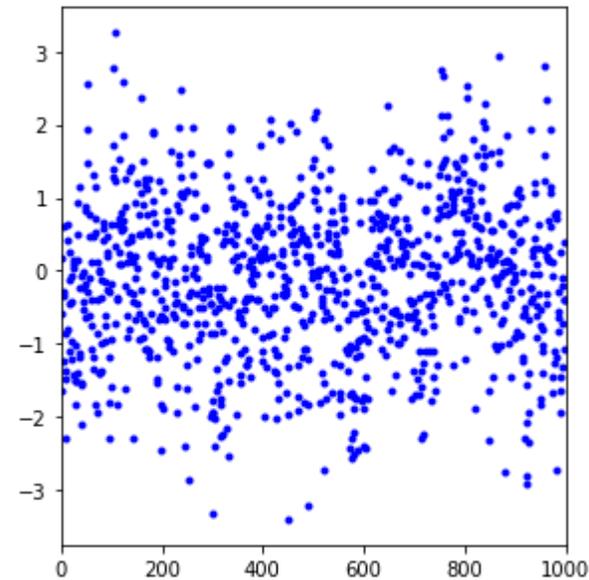
# Champs gaussiens aléatoires : exemples

$$C_{ij} = \frac{1}{(1 + |i - j|/2)^2}$$

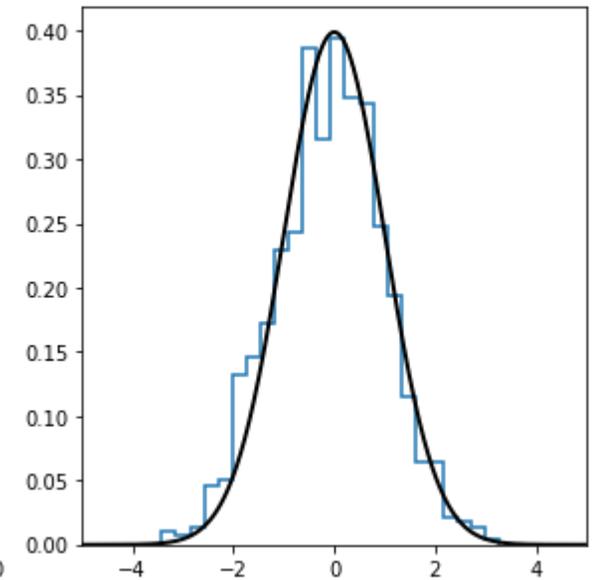
matrice de covariance



champ gaussien



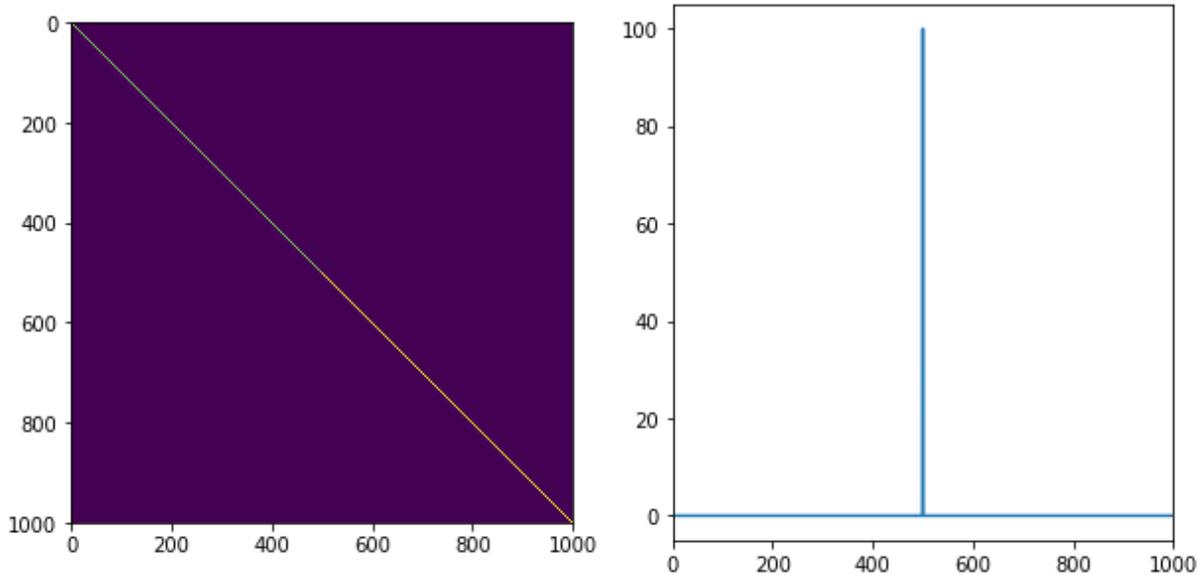
histogramme



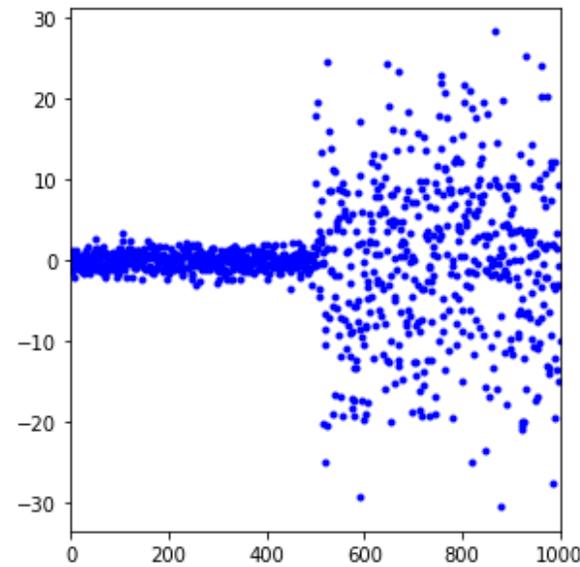
## Champs gaussiens aléatoires : exemples

$$C_{ii} = \begin{cases} 1 & \text{if } i < N/2 \\ 100 & \text{otherwise} \end{cases}$$
$$C_{ij} = 0 \quad \text{for } i \neq j$$

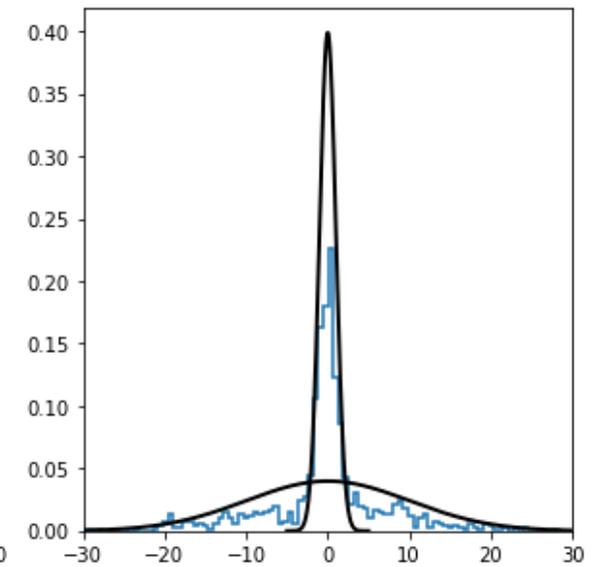
matrice de covariance



champ gaussien



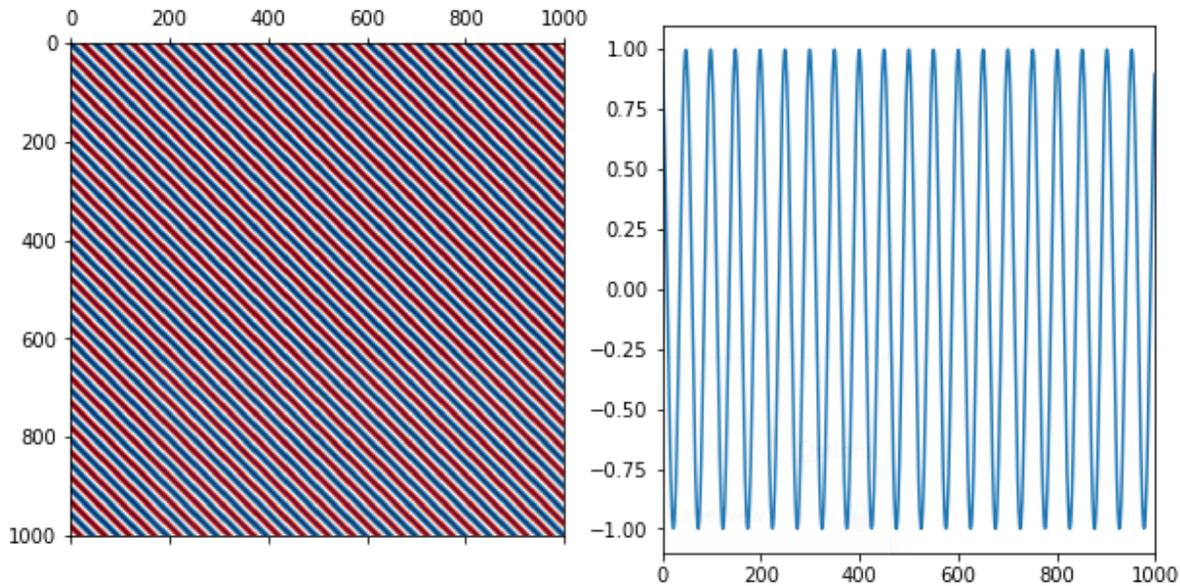
histogramme



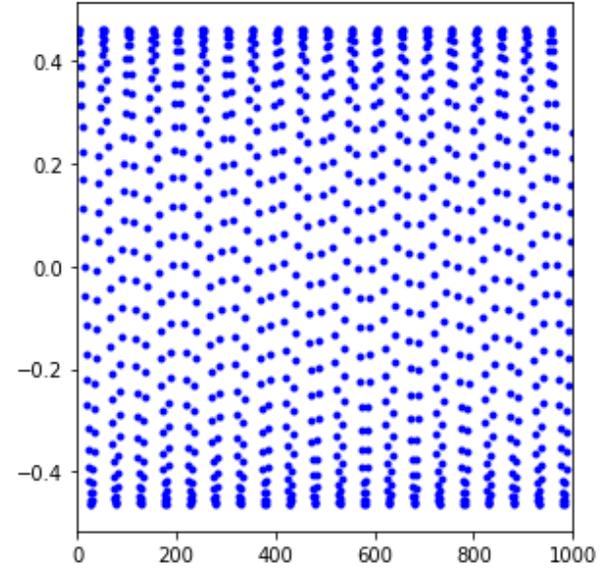
# Champs gaussiens aléatoires : exemples

$$C_{ij} = \cos\left(\frac{i-j}{8}\right)$$

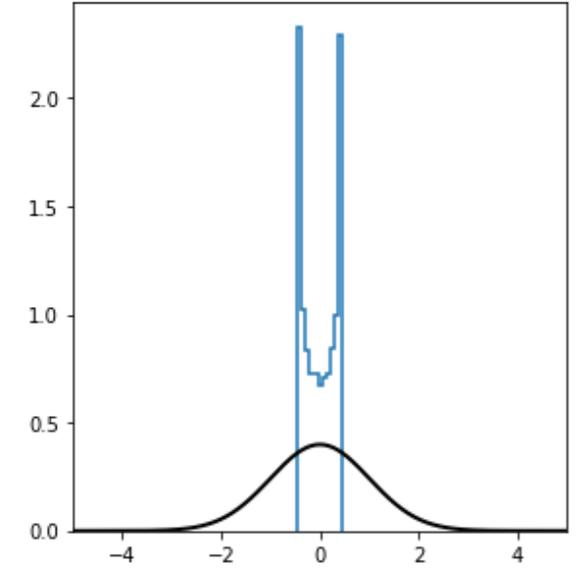
matrice de covariance



champ gaussien



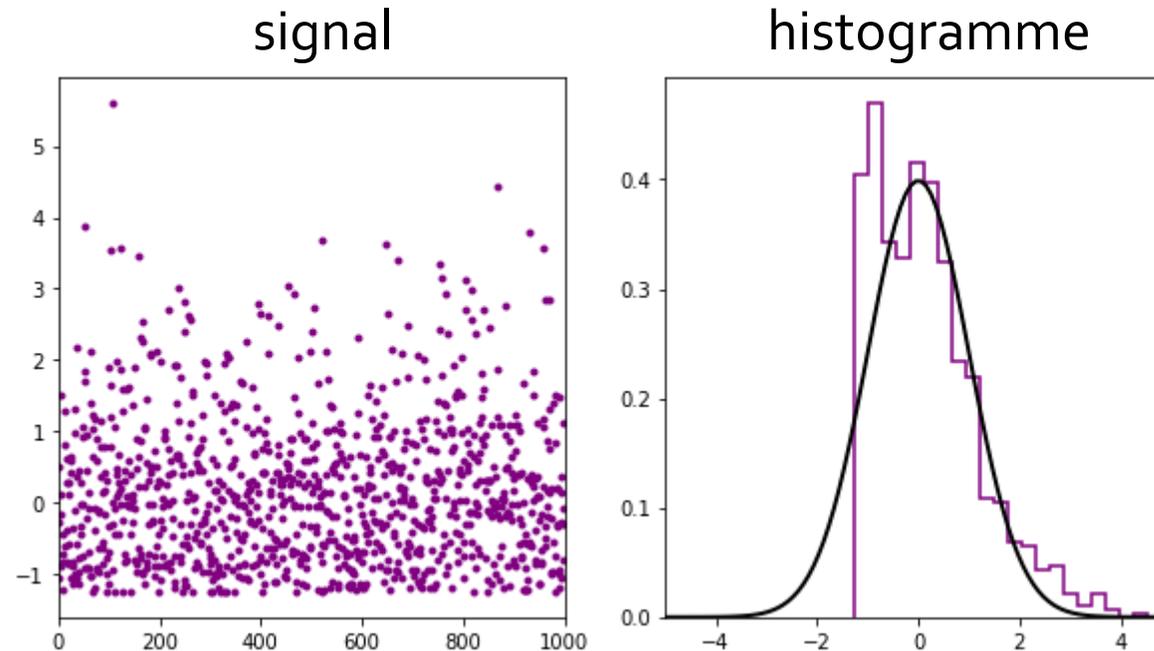
histogramme



Les histogrammes de champs aléatoires gaussiens ne sont pas toujours gaussiens !

## Exemple de signal non-Gaussien

- $s = \Phi + f_{\text{NL}}\Phi^2$  où  $\Phi$  est un champ gaussien aléatoire.
- En cosmologie, cela s'appelle « non-gaussianité de type local ».
- Les non-gaussianités primordiales renseignent sur la physique de l'inflation cosmologique.



- La pdf à un point possède une asymétrie (*skewness*).

## Probabilités marginales et conditionnelles de champs gaussiens aléatoires

- On travaille avec deux variables qui constituent **conjointement** un champ gaussien  $\begin{pmatrix} x \\ y \end{pmatrix}$ .
- Alors  $x$  et  $y$  sont des champs gaussiens dont les **probabilités marginales** ont pour

$$\text{Moyenne : } \mu = \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix} \quad \text{Matrice de covariance : } C = \begin{pmatrix} C_{xx} & C_{xy} \\ C_{yx} & C_{yy} \end{pmatrix}$$

(Les moyennes et covariances marginales sont simplement les parties correspondantes de la moyenne et de la covariance jointe.)

- Les **probabilités conditionnelles** sont également gaussiennes avec :

$$\begin{aligned} \text{Moyenne : } \quad \mu_{x|y} &= \mu_x + C_{xy}C_{yy}^{-1}(y - \mu_y) \\ \text{Matrice de covariance : } \quad C_{x|y} &= C_{xx} - C_{xy}C_{yy}^{-1}C_{yx} \end{aligned}$$

# Traitement bayésien du signal et filtrage de Wiener

Notebook : [https://github.com/florent-](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising.ipynb)

[leclercq/Bayes\\_InfoTheory/blob/master/WienerFilter\\_denoising.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising.ipynb)

Notebook : [https://github.com/florent-](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising_CMB.ipynb)

[leclercq/Bayes\\_InfoTheory/blob/master/WienerFilter\\_denoising\\_CMB.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising_CMB.ipynb)

[nb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_denoising_CMB.ipynb)

Notebook : [https://github.com/florent-](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_deblending.ipynb)

[leclercq/Bayes\\_InfoTheory/blob/master/WienerFilter\\_deblending.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/WienerFilter_deblending.ipynb)

## Le filtrage de Wiener : débruitage bayésien d'un signal

- Modèle de données :  $d = s + n$  où  $\begin{pmatrix} s \\ d \end{pmatrix}$  est un champ gaussien aléatoire (conjointement).

- Solution :

$$\mu_{s|d} = \mu_s + C_{sd}C_{dd}^{-1} (d - \mu_d)$$

$$C_{s|d} = C_{ss} - C_{sd}C_{dd}^{-1}C_{ds}$$

- Notations :  $C_{ss} \equiv S$  et  $C_{nn} \equiv N$ .
- Hypothèses supplémentaires :  $C_{sn} = C_{ns} = 0$ . Alors :

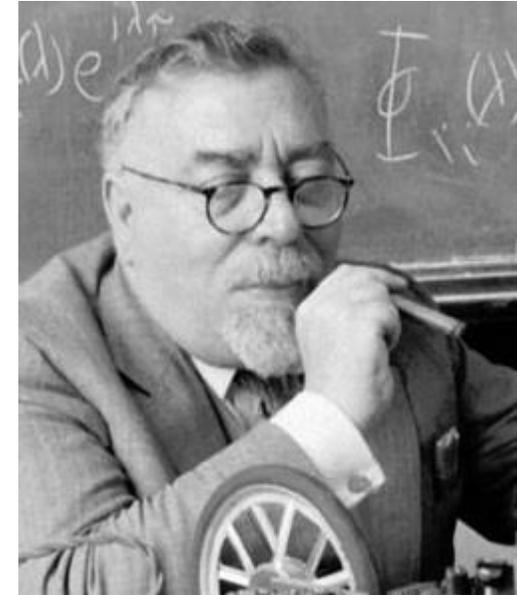
$$C_{dd} = S + N$$

$$C_{sd} = C_{ss} + C_{sn} = C_{ss} = S$$

- Expressions finales :

$$\mu_{s|d} = \mu_s + S(S + N)^{-1}(d - \mu_d) = \mu_s + (S^{-1} + N^{-1})^{-1}N^{-1}(d - \mu_d)$$

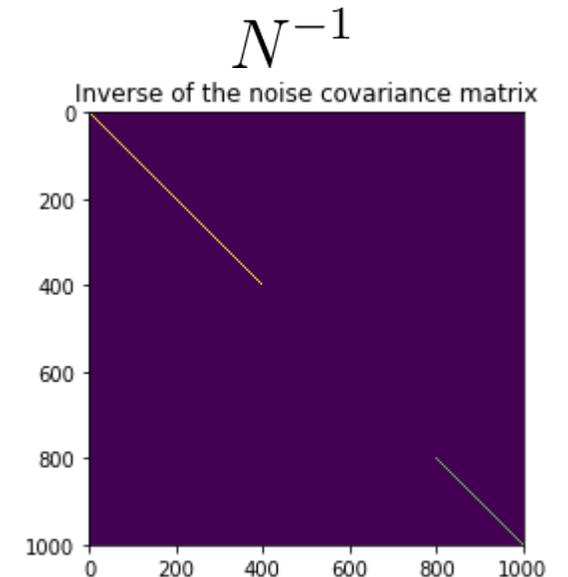
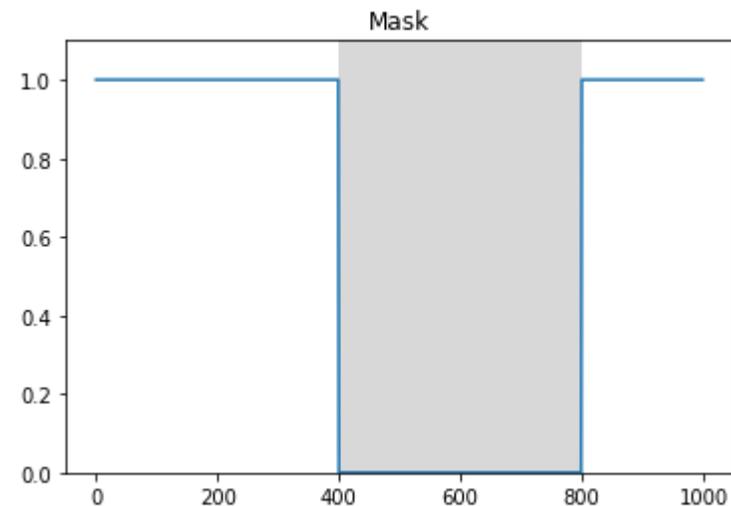
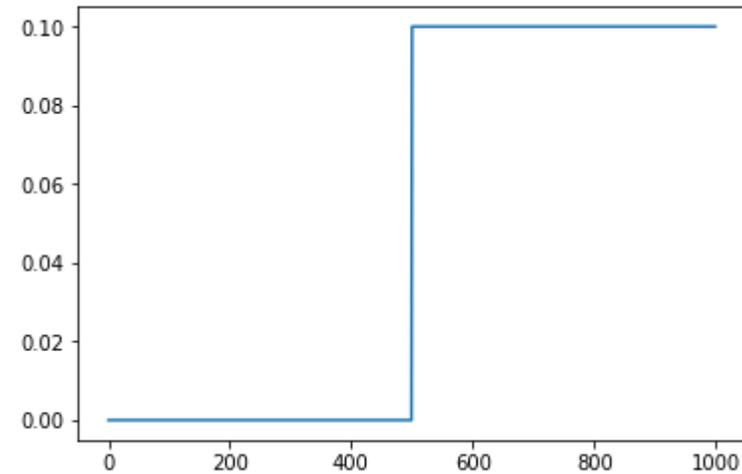
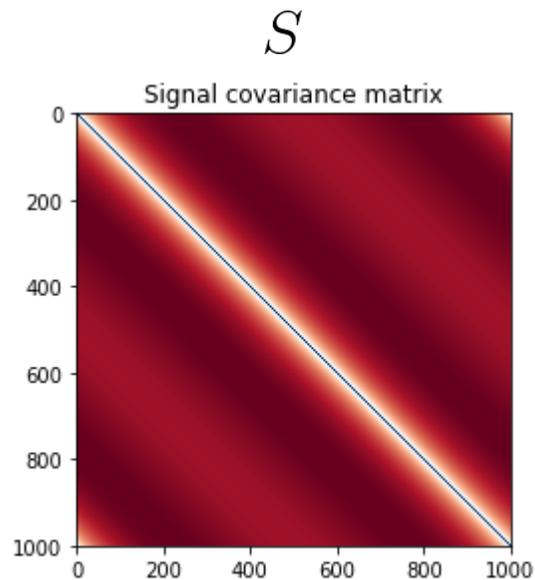
$$C_{s|d} = S - S(S + N)^{-1}S = (S^{-1} + N^{-1})^{-1}$$



Norbert Wiener  
(1894-1964)

## Exemple de débruitage par filtrage de Wiener

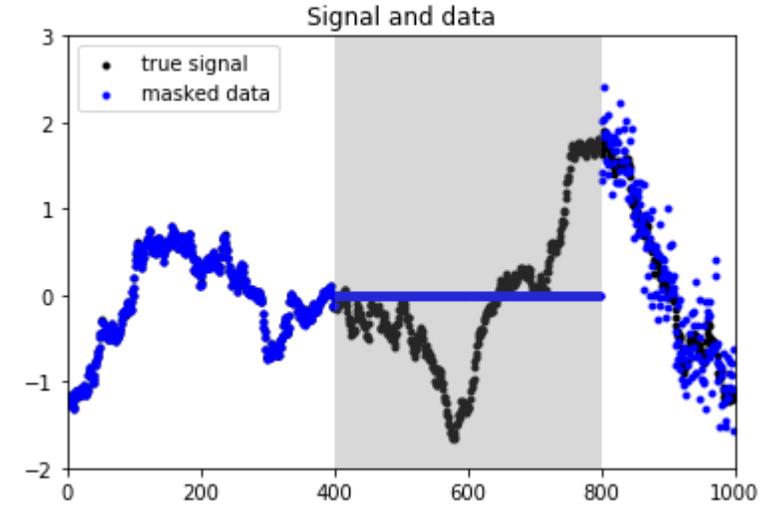
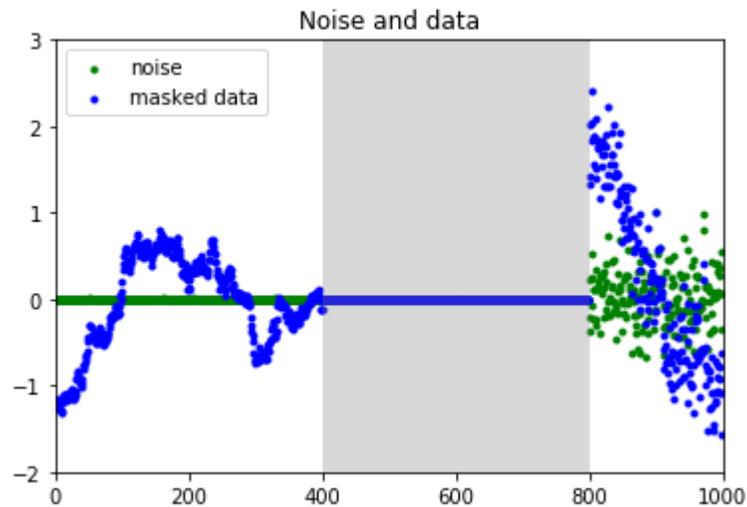
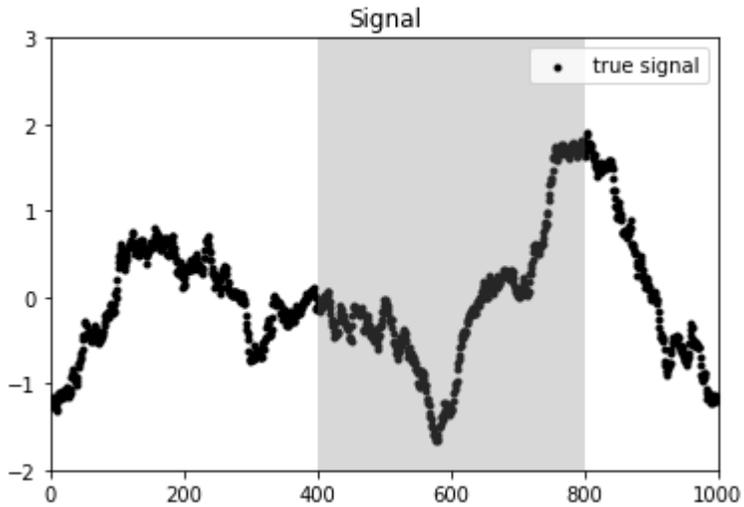
- Hypothèses concernant la matrice de covariance du signal et du bruit



# Exemple de débruitage par filtrage de Wiener

- Génération d'un signal et de données synthétiques

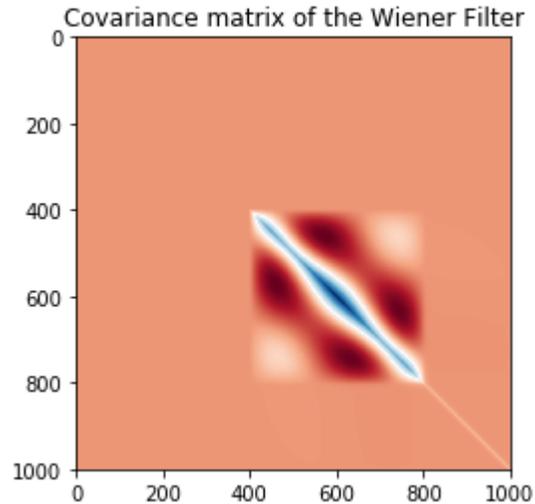
$$d = s + n$$



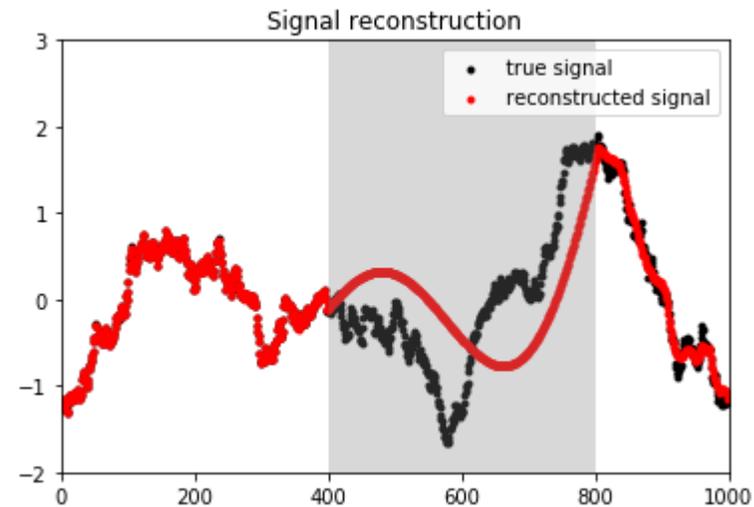
## Exemple de débruitage par filtrage de Wiener

- Filtrage de Wiener
- La moyenne de la reconstruction correspond à l'estimateur maximum a posteriori.

$$C_{s|d} = (S^{-1} + N^{-1})^{-1}$$



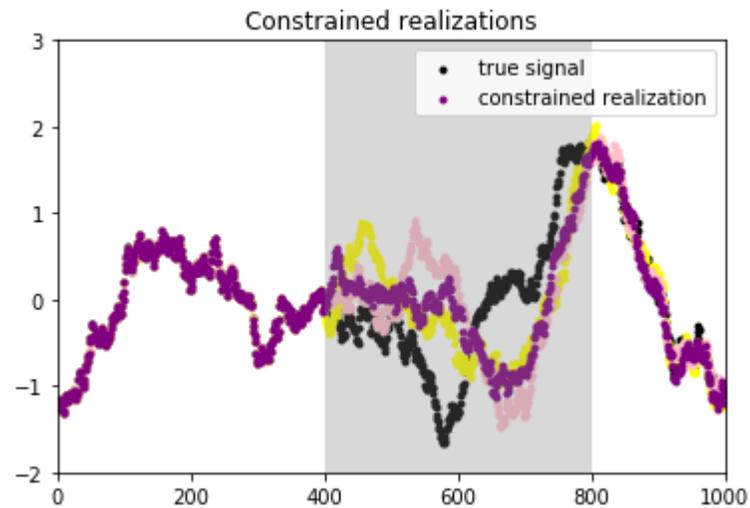
$$\mu_{s|d} = \mu_s + C_{s|d}N^{-1}(d - \mu_d)$$



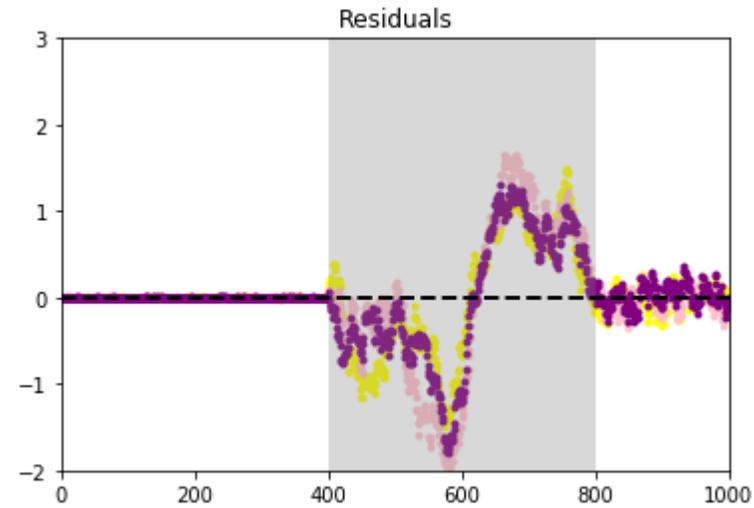
## Exemple de débruitage par filtrage de Wiener

- Simulations contraintes du signal débruité

$$s_{\text{sim}} = \mu_{s|d} + \sqrt{C_{s|d}} \xi$$



$$s_{\text{sim}} - s_{\text{true}}$$



## Le filtrage de Wiener : séparation bayésienne de composantes mélangées et bruitées

- Modèle de données : 
$$d = \begin{pmatrix} x_1 \\ x_2 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} n_1 \\ n_2 \\ n_3 \end{pmatrix}$$

- Hypothèses : 
$$C_{x_1d} = \begin{pmatrix} C_{x_1x_1} & C_{x_1x_1} & 0 \\ 0 & C_{x_2x_2} & C_{x_2x_2} \end{pmatrix} \quad C_{nn} = \begin{pmatrix} C_{n_1n_1} & 0 & 0 \\ 0 & C_{n_2n_2} & 0 \\ 0 & 0 & C_{n_3n_3} \end{pmatrix}$$

$$C_{dd} = \begin{pmatrix} C_{x_1x_1} + C_{n_1n_1} & C_{x_1x_1} & 0 \\ C_{x_1x_1} & C_{x_1x_1} + C_{x_2x_2} + C_{n_2n_2} & C_{x_2x_2} \\ 0 & C_{x_2x_2} & C_{x_2x_2} + C_{n_3n_3} \end{pmatrix}$$

- Solution :

$$\mu_{x_1|d} = C_{x_1d}C_{dd}^{-1}d$$

$$\mu_{x_2|d} = C_{x_2d}C_{dd}^{-1}d$$

$$C_{x_1|d} = C_{x_1x_1} - C_{x_1d}C_{dd}^{-1}C_{dx_1}$$

$$C_{x_2|d} = C_{x_2x_2} - C_{x_2d}C_{dd}^{-1}C_{dx_2}$$

# Méthodes Monte Carlo et chaînes de Markov

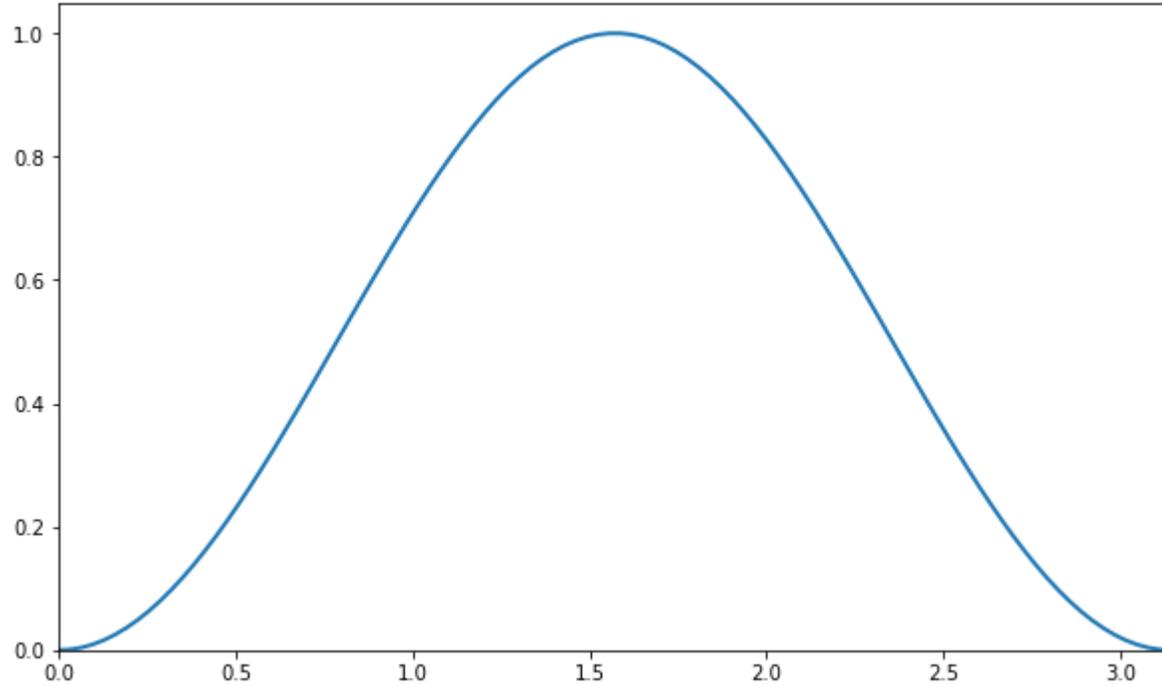
Notebook : [https://github.com/florent-leclercq/Bayes\\_InfoTheory/blob/master/Sampling\\_Importance\\_Rejection.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/Sampling_Importance_Rejection.ipynb)

Notebook : [https://github.com/florent-leclercq/Bayes\\_InfoTheory/blob/master/MCMC\\_MH.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/MCMC_MH.ipynb)

Notebook : [https://github.com/florent-leclercq/Bayes\\_InfoTheory/blob/master/Supernova\\_MCMC\\_HMC.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/Supernova_MCMC_HMC.ipynb)

# Intégration Monte Carlo : échantillonnage standard

Target pdf



$$I = \int_a^b f(x) dx$$

```
In [6]: trueI=quad(target_pdf,a,b)[0]
trueI
```

```
Out[6]: 1.5707963267948966
```

Standard Monte Carlo integration

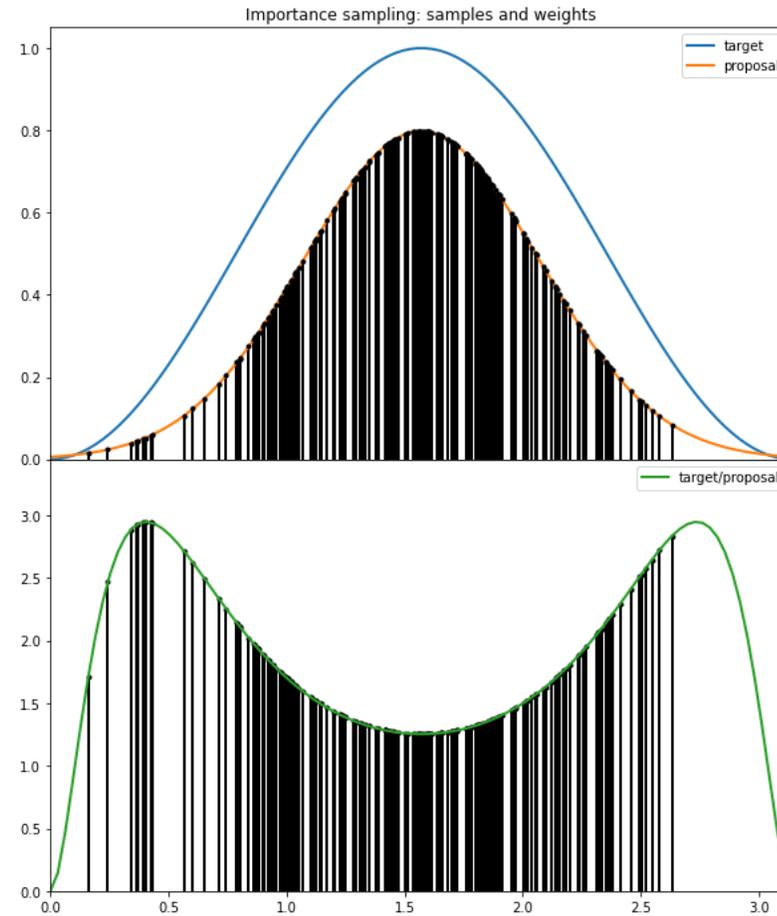
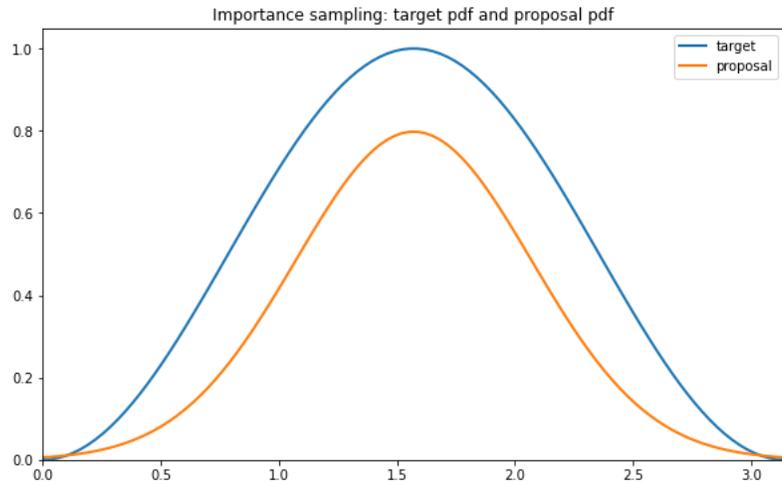


$$I \approx \frac{b-a}{N_{\text{samples}}} \sum_{i=1}^{N_{\text{samples}}} x_i$$

```
In [9]: StandardMonteCarloI=np.sum(samples)*(b-a)/Nsamp
StandardMonteCarloI
```

```
Out[9]: 1.660402945484072
```

# Intégration Monte Carlo : échantillonnage préférentiel (*importance sampling*)



$$I \approx \frac{\sum_{i=1}^{N_{\text{samples}}} x_i w_i}{\sum_{i=1}^{N_{\text{samples}}} w_i}$$

```
In [14]: ImportanceI=np.average(samples,weights=weights)
ImportanceI
```

```
Out[14]: 1.5100957559182684
```

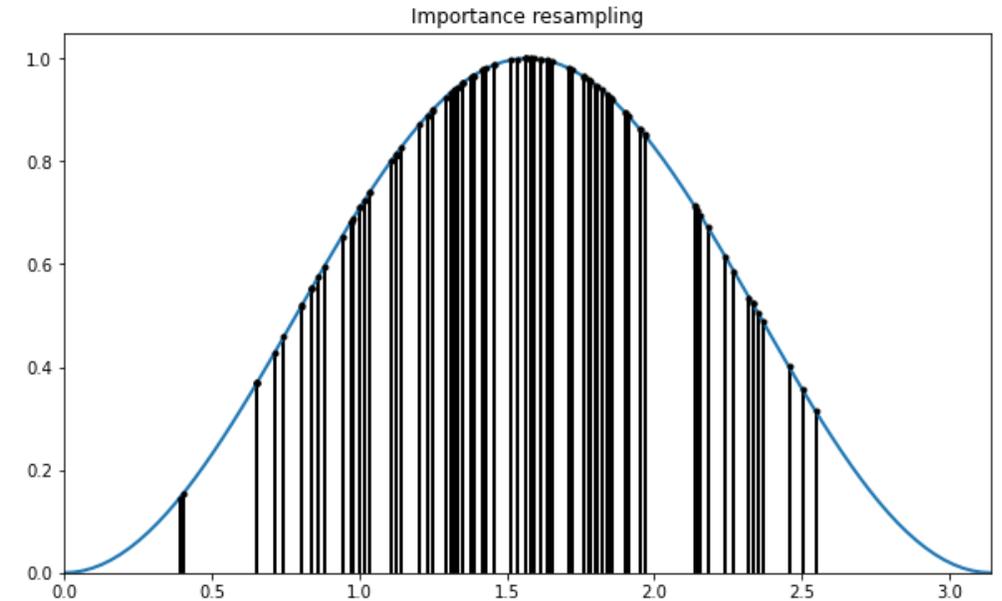
## Importance resampling

A problem with importance sampling is the situation in which all but one of the weights are close to zero. To avoid with situation, we can do **importance resampling**. We draw  $N_{\text{resamp}}$  new samples from the current sample set with probabilities proportional to their weights. We replace the current samples with this new set, and the the current weights by  $1/N_{\text{resamp}}$  (drawing according to the importance weight replaces likelihoods by frequencies).

```
In [15]: Nresamp=100
normalizedweights=weights/np.sum(weights)
resamples=np.random.choice(samples, size=Nresamp, replace=True, p=normalizedweights)
reweights=1./Nresamp*np.ones(Nresamp)
```

Weights are then updated given their likelihood, as in the previous importance sampling step.

```
In [16]: reweights*=target_pdf(resamples)/(1./Nresamp)
```

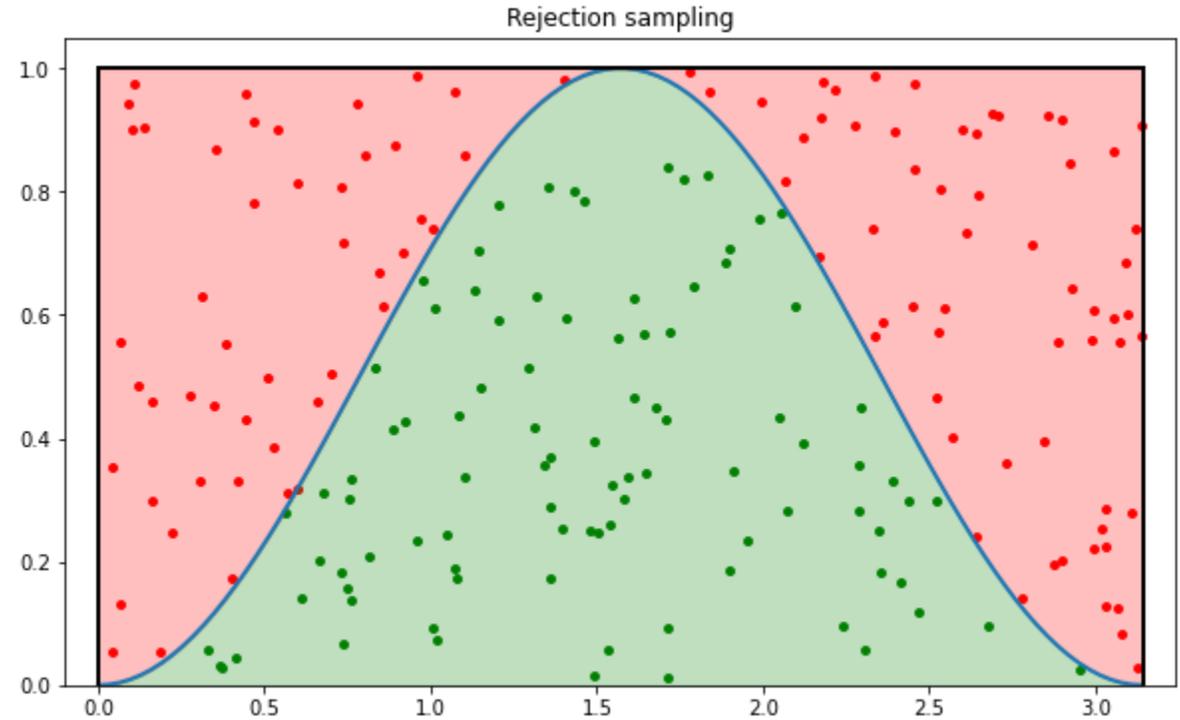
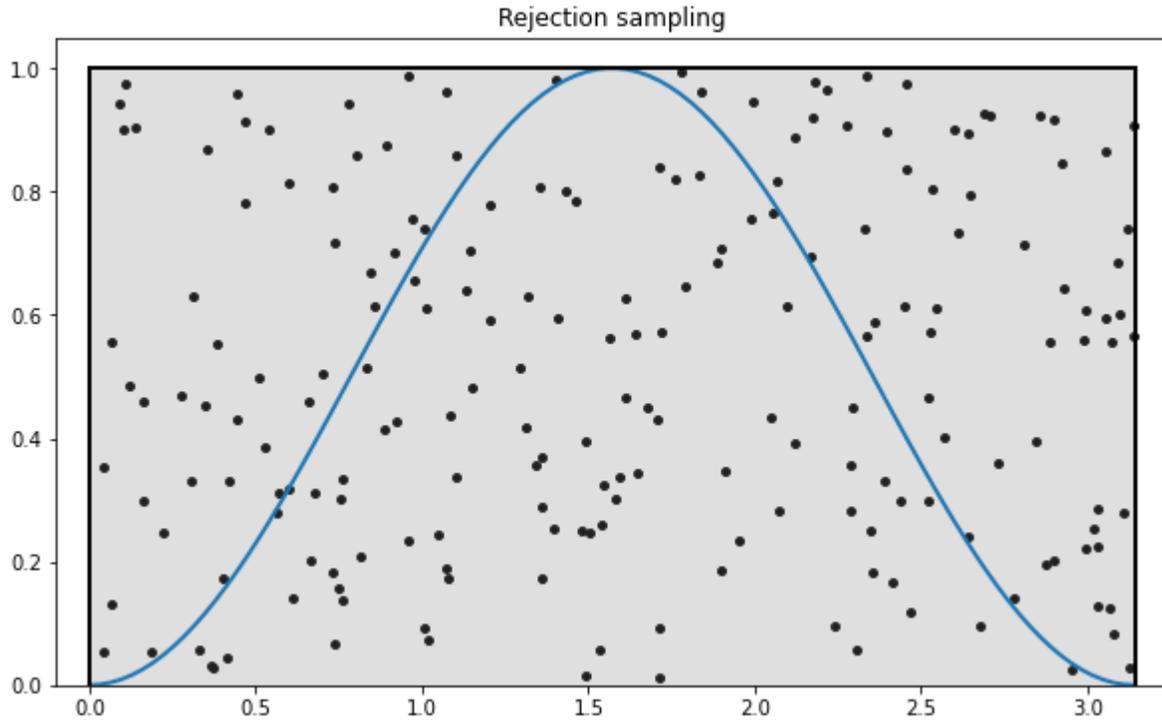


$$I \approx \frac{\sum_{i=1}^{N_{\text{samples}}} x_i w_i}{\sum_{i=1}^{N_{\text{samples}}} w_i}$$

```
In [18]: ImportanceReI=np.average(resamples,weights=reweights)
ImportanceReI
```

```
Out[18]: 1.531675529534044
```

# Intégration Monte Carlo : échantillonnage par rejet (*rejection sampling*)



```
In [23]: fraction=float(len(accepted_samples))/Nsamp  
fraction
```

```
Out[23]: 0.465
```

```
In [24]: fraction=float(len(accepted_samples))/Nsamp  
RejectionI=fraction*upperbound*(b-a)  
RejectionI
```

```
Out[24]: 1.460840583919254
```

# Chaînes de Markov Monte Carlo (Markov Chain Monte Carlo – MCMC)

- Propriété de Markov :
  - L'information utile pour la prédiction du futur est entièrement contenue dans l'état présent du processus et n'est pas dépendante des états antérieurs (le système n'a pas de « mémoire »).
  - Mathématiquement : la distribution conditionnelle de probabilité des états futurs, étant donnés les états passés et l'état présent, ne dépend que de l'état présent et non pas des états passés.
- Algorithme de Metropolis-Hastings :

```
begin
  initialise  $\mathbf{x}_{(0)}$ ;
  for  $i = 1$  to  $n$  do
     $\mathbf{x}^* \leftarrow q(\mathbf{x}^*|\mathbf{x})$  (proposal distribution);
     $\alpha \leftarrow \mathcal{U}(0, 1)$  (uniform distribution);
    if  $\alpha < \min [1, r(\mathbf{x}, \mathbf{x}^*)]$  then
      |  $\mathbf{x}_{(i)} = \mathbf{x}^*$ ;
    else
      |  $\mathbf{x}_{(i)} = \mathbf{x}_{(i-1)}$ ;
    end
  end
  return  $(\mathbf{x}_{(0)}, \dots, \mathbf{x}_{(n)})$ ;
end
```

$$\text{Cas général : } r(\mathbf{x}, \mathbf{x}^*) \equiv \frac{p(\mathbf{x}^*) q(\mathbf{x}|\mathbf{x}^*)}{p(\mathbf{x}) q(\mathbf{x}^*|\mathbf{x})} \quad (\text{ratio de Hastings})$$

$$\text{Cas particulier : } r(\mathbf{x}, \mathbf{x}^*) = \frac{p(\mathbf{x}^*)}{p(\mathbf{x})} \quad (\text{Metropolis update})$$

dans le cas d'un *proposal* symétrique :  $q(\mathbf{x}^*|\mathbf{x}) = q(\mathbf{x}|\mathbf{x}^*)$

# Algorithme de Metropolis-Hastings : implémentation

```
begin
  initialise  $\mathbf{x}_{(0)}$ ;
  for  $i = 1$  to  $n$  do
     $\mathbf{x}^* \leftarrow q(\mathbf{x}^*|\mathbf{x})$  (proposal distribution);
     $\alpha \leftarrow \mathcal{U}(0, 1)$  (uniform distribution);
    if  $\alpha < \min[1, r(\mathbf{x}, \mathbf{x}^*)]$  then
      |  $\mathbf{x}_{(i)} = \mathbf{x}^*$ ;
    else
      |  $\mathbf{x}_{(i)} = \mathbf{x}_{(i-1)}$ ;
    end
  end
  return  $(\mathbf{x}_{(0)}, \dots, \mathbf{x}_{(n)})$ ;
end
```

```
def MH_sampler(Ntries, theta_start, Ntrials, Nsuccesses, lh, prior, proposal_sigma):
    Naccepted=0
    samples=np.zeros(Ntries+1)
    samples[0]=theta_start
    theta=theta_start
    for i in range(Ntries):
        theta_p = theta + proposal_pdf(proposal_sigma).rvs()
        # the Gaussian proposal pdf satisfies the detailed balance equation, so the
        # acceptance ratio simplifies to the Metropolis ratio
        a = min(1, target_pdf(theta_p, Ntrials, Nsuccesses, lh, prior)/target_pdf(theta, Ntrials, Nsuccesses, lh, prior))
        u = np.random.uniform()
        if u < a:
            Naccepted+=1
            theta=theta_p
        samples[i+1] = theta
    return Naccepted, samples
```

Cas général :  $r(\mathbf{x}, \mathbf{x}^*) \equiv \frac{p(\mathbf{x}^*) q(\mathbf{x}|\mathbf{x}^*)}{p(\mathbf{x}) q(\mathbf{x}^*|\mathbf{x})}$  (ratio de Hastings)

Cas particulier :  $r(\mathbf{x}, \mathbf{x}^*) = \frac{p(\mathbf{x}^*)}{p(\mathbf{x})}$  (Metropolis update)

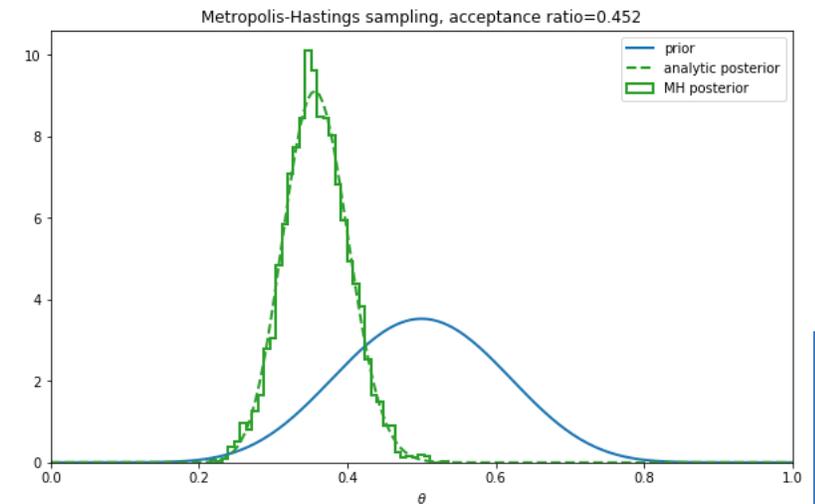
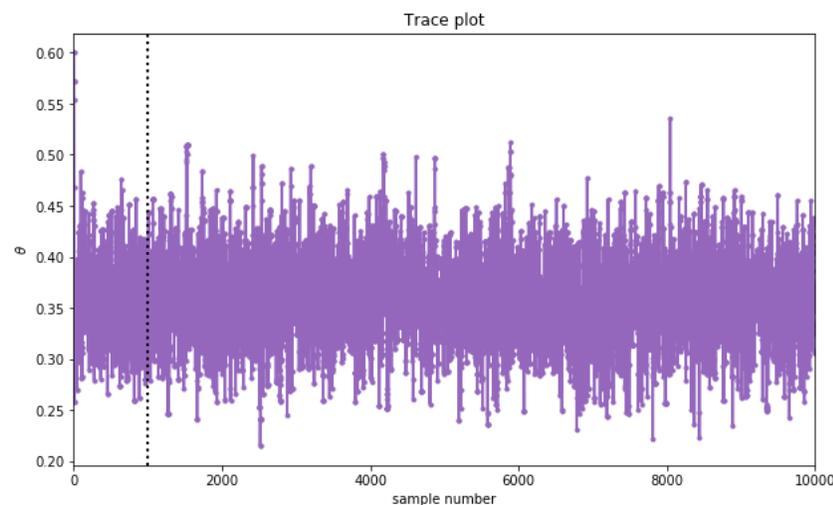
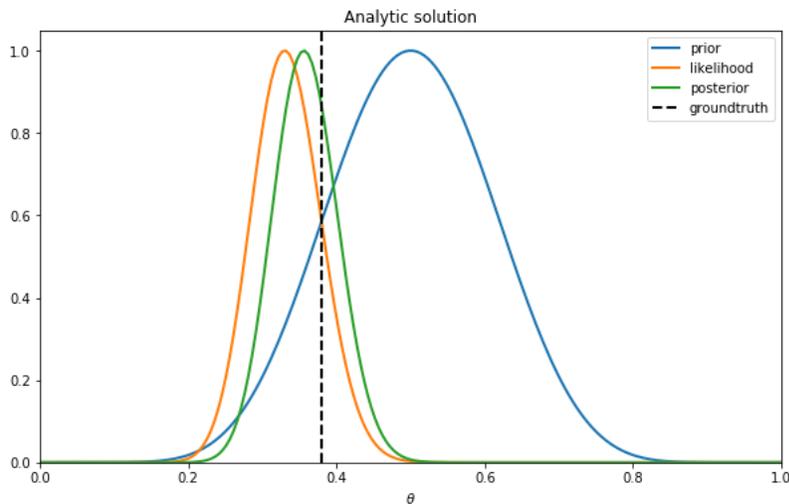
dans le cas d'un *proposal* symétrique :  $q(\mathbf{x}^*|\mathbf{x}) = q(\mathbf{x}|\mathbf{x}^*)$

## Exemple de MCMC

- Problème considéré : une expérience de Bernoulli ( $N_{\text{trials}}$  expériences indépendantes avec une probabilité de succès  $\theta$ ).
- La vraisemblance pour ce problème est la distribution binomiale :

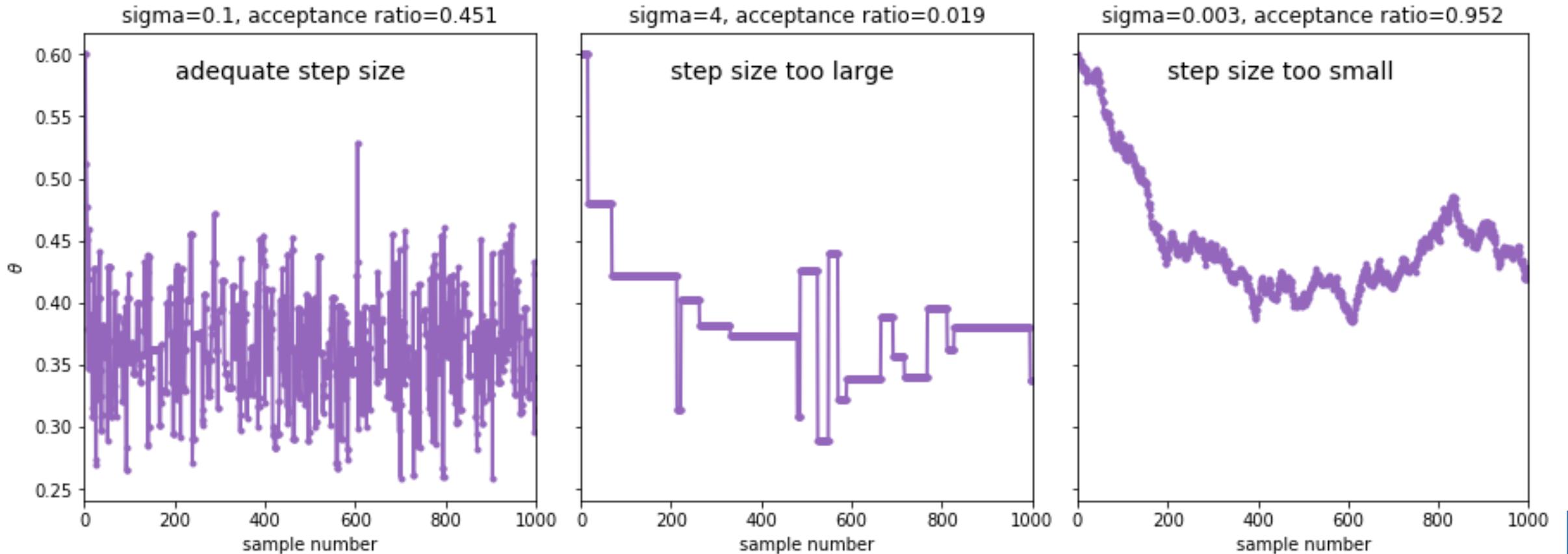
$$p(N_{\text{successes}}, N_{\text{trials}}, \theta) = \binom{N_{\text{trials}}}{N_{\text{successes}}} \theta^{N_{\text{successes}}} (1 - \theta)^{N_{\text{trials}} - N_{\text{successes}}}$$

- Résultat analytique : la distribution beta est une famille de **priors conjugués**, c'est-à-dire : si le prior est  $\mathcal{B}(\alpha, \beta)$ , alors le posterior est  $\mathcal{B}(\alpha', \beta')$  avec  $\alpha' = \alpha + N_{\text{successes}}$   
 $\beta' = \beta + N_{\text{trials}} - N_{\text{successes}}$



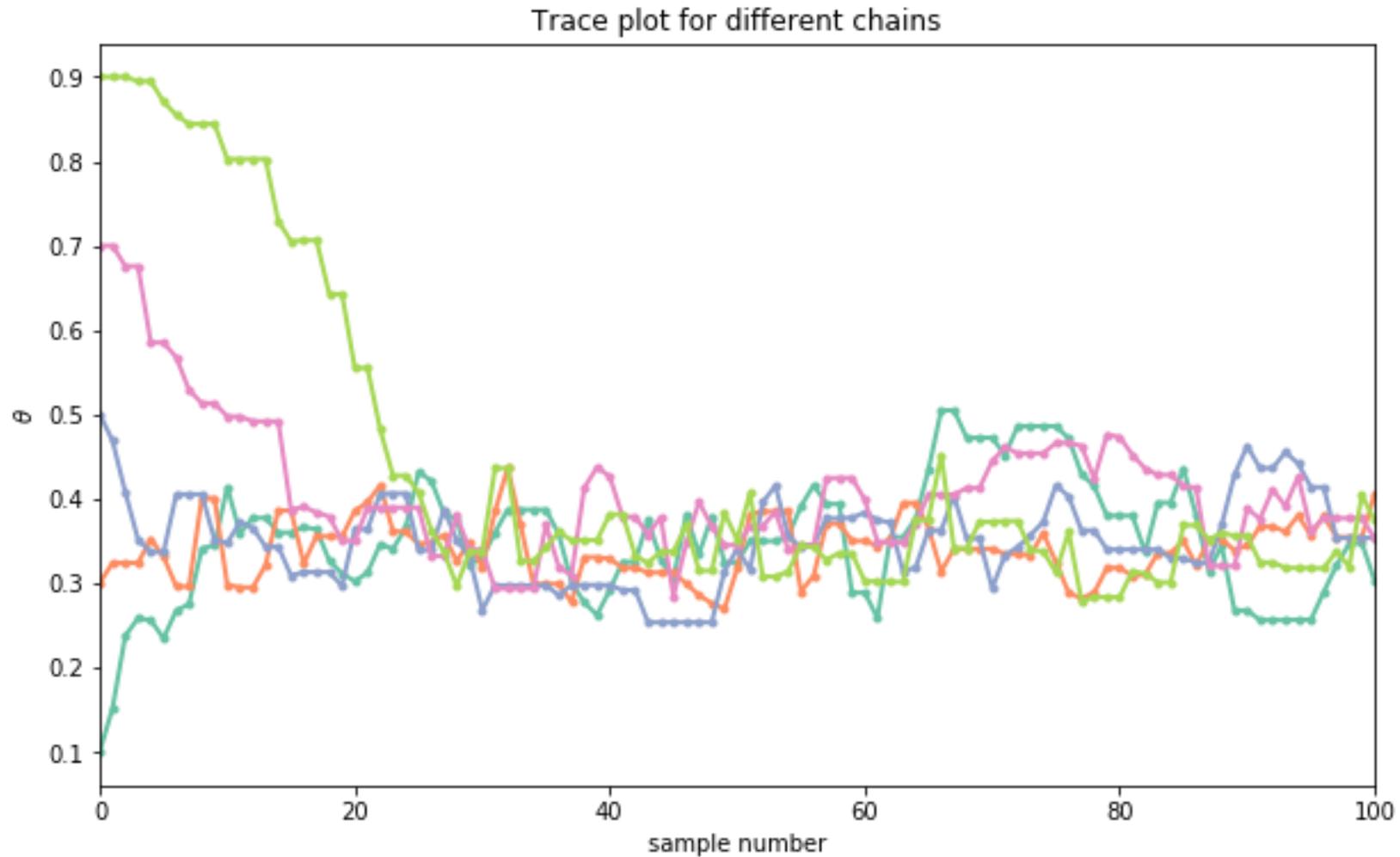
## Diagnostics des chaînes de Markov : les *trace plots*

- L'ajustement du *proposal* :



## Diagnostics des chaînes de Markov : les *trace plots*

- Plusieurs chaînes indépendantes, différents points de départ :



## Diagnostiques des chaînes de Markov : convergence – le test de Gelman-Rubin

- Paramètres:

- $m$  : nombre de chaînes
- $n$  : longueur des chaînes

- Définitions :

- Variance inter-chaînes ("*between*" chains variance) :  $B \equiv \frac{n}{m-1} \sum_{j=1}^m (\bar{\psi}_{.j} - \bar{\psi}_{..})^2$

$$\bar{\psi}_{.j} = \frac{1}{n} \sum_{i=1}^n \psi_{ij}$$

$$\bar{\psi}_{..} = \frac{1}{m} \sum_{j=1}^m \bar{\psi}_{.j}$$

- Variance intra-chaînes ("*within*" chains variance) :  $W \equiv \frac{1}{m} \sum_{j=1}^m s_j^2$   $s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\psi_{ij} - \bar{\psi}_{.j})^2$

- Estimateurs (de la variance du posterior marginal de chacun des paramètres) :

- $\widehat{\text{var}}^- \equiv W$  sous-estime la variance
- $\widehat{\text{var}}^+ \equiv \frac{n}{n-1} W + \frac{1}{n} B$  sur-estime la variance

- Test de Gelman-Rubin :

- Facteur de réduction d'échelle potentiel (*Potential scale reduction factor* – PSRF) :  $\widehat{R} \equiv \sqrt{\frac{\widehat{\text{var}}^+}{\widehat{\text{var}}^-}}$
- Test :  $\widehat{R} \rightarrow 1$  quand  $n \rightarrow \infty$ . On vise typiquement  $\widehat{R} - 1 \lesssim 10^{-2}$ .

# Priors d'ignorance et le principe de l'entropie maximale

Notebook : [https://github.com/florent-](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/LighthouseProblem.ipynb)

[leclercq/Bayes\\_InfoTheory/blob/master/LighthouseProblem.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/LighthouseProblem.ipynb)

Notebook : [https://github.com/florent-](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/MaximumEntropy.ipynb)

[leclercq/Bayes\\_InfoTheory/blob/master/MaximumEntropy.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/MaximumEntropy.ipynb)

- **Priors d'ignorance** : l'idée est d'imposer l'invariance de l'état de connaissance selon une transformation :

$$p(T(x))dT(x) = p(x)dx$$

- Cas le plus simple : symétrie sous l'échange de deux modèles  $\mathcal{M}_1$  et  $\mathcal{M}_2$  :

$$\begin{aligned} p(\mathcal{M}_1) &= p(\mathcal{M}_2) \\ p(\mathcal{M}_1) + p(\mathcal{M}_2) &= 1 \end{aligned} \quad \Rightarrow \quad p(\mathcal{M}_1) = p(\mathcal{M}_2) = \frac{1}{2} \quad \text{Symétrie } \mathbb{Z}_2$$

- Incertitude maximale sur un « paramètre de localisation » de la distribution :  $T(x) = x + s \quad \forall s$

$$\begin{aligned} dT &= dx \\ p(x) &= p(x + s) \quad \forall s \end{aligned} \quad \Rightarrow \quad p(x) = C \quad \text{Prior plat} \quad \text{Symétrie } U(1)$$

- Incertitude maximale sur un « paramètre d'échelle » de la distribution :  $T(x) = ax \quad \forall a$

$$\begin{aligned} dT &= a dx \\ p(x) &= a p(ax) \quad \forall a \end{aligned} \quad \Rightarrow \quad p(x) = C/a \quad \text{Prior de Jeffreys} \quad \text{Symétrie } U(1)$$

- Cas général : spécifier un **groupe de transformations** et résoudre une **équation fonctionnelle**.

# Le problème du phare (the lighthouse problem)



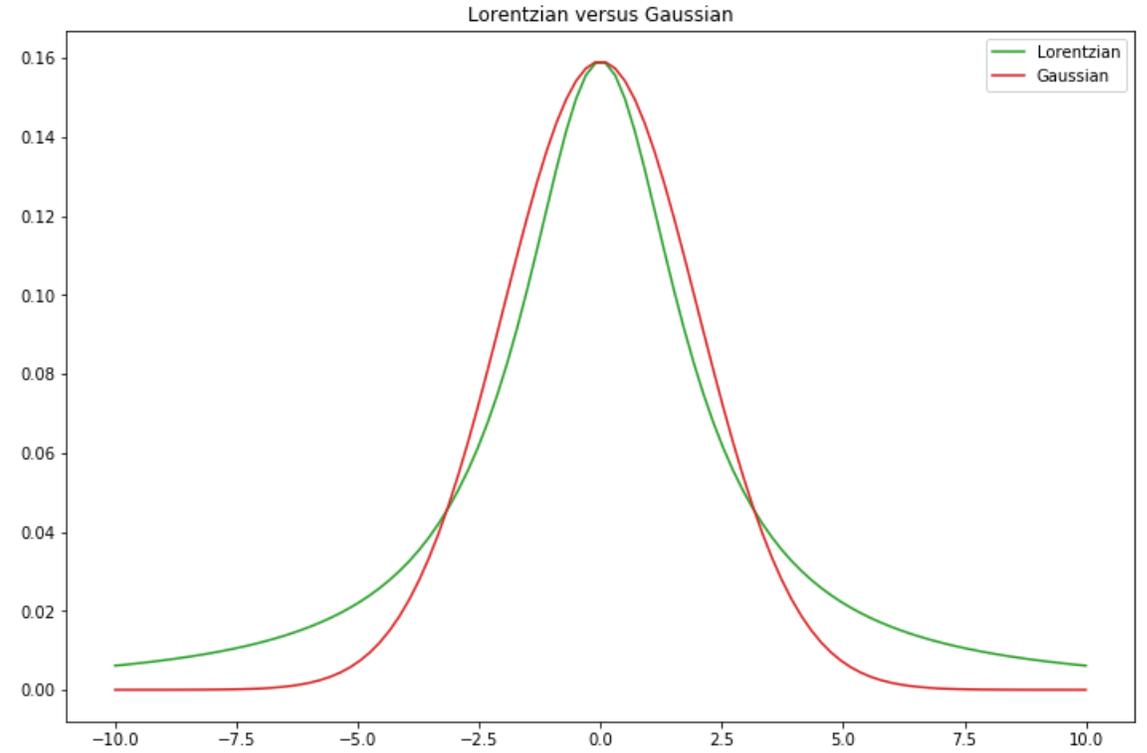
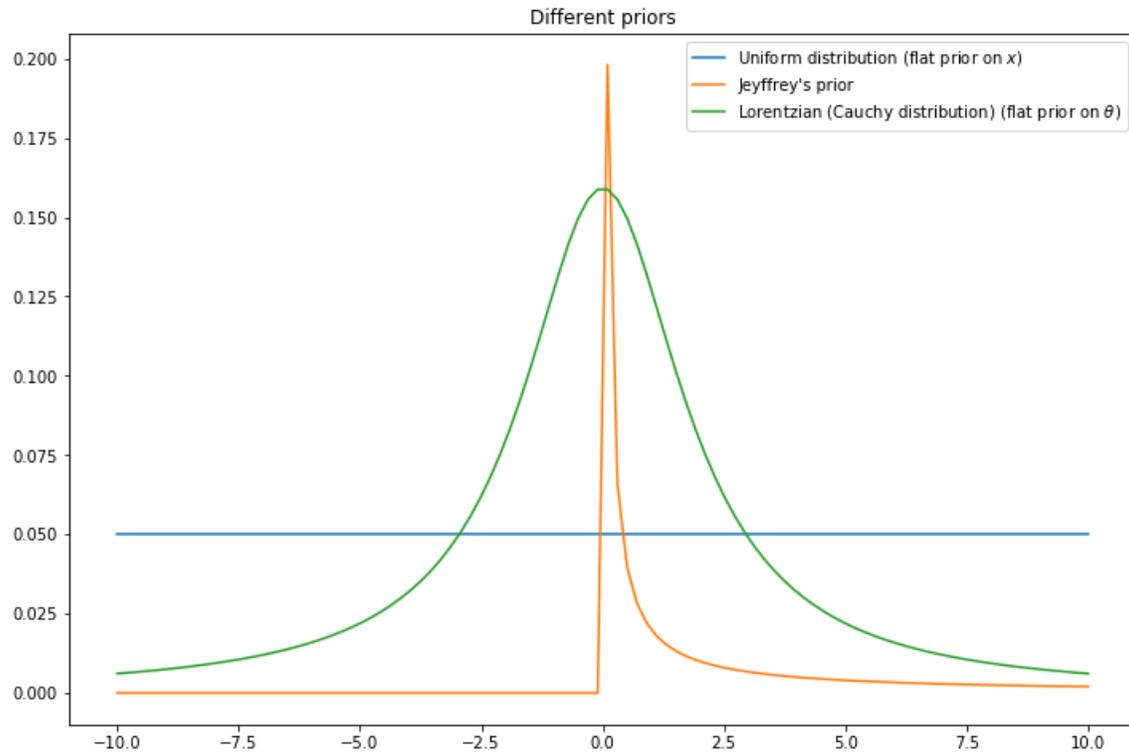
$$x = d \tan \theta$$

$$dx = d(1 + \tan^2 \theta)d\theta \quad \text{Si } p(\theta) = C \quad \Rightarrow$$

$$p(x) \propto \frac{d}{d^2 + x^2}$$

$$p(x)dx = p(\theta)d\theta$$

## Lorentzienne/Distribution de Cauchy



- Un état d'ignorance maximale pour une variable n'est généralement pas la même chose qu'un état d'ignorance maximale pour une fonction non-linéaire de cette variable.

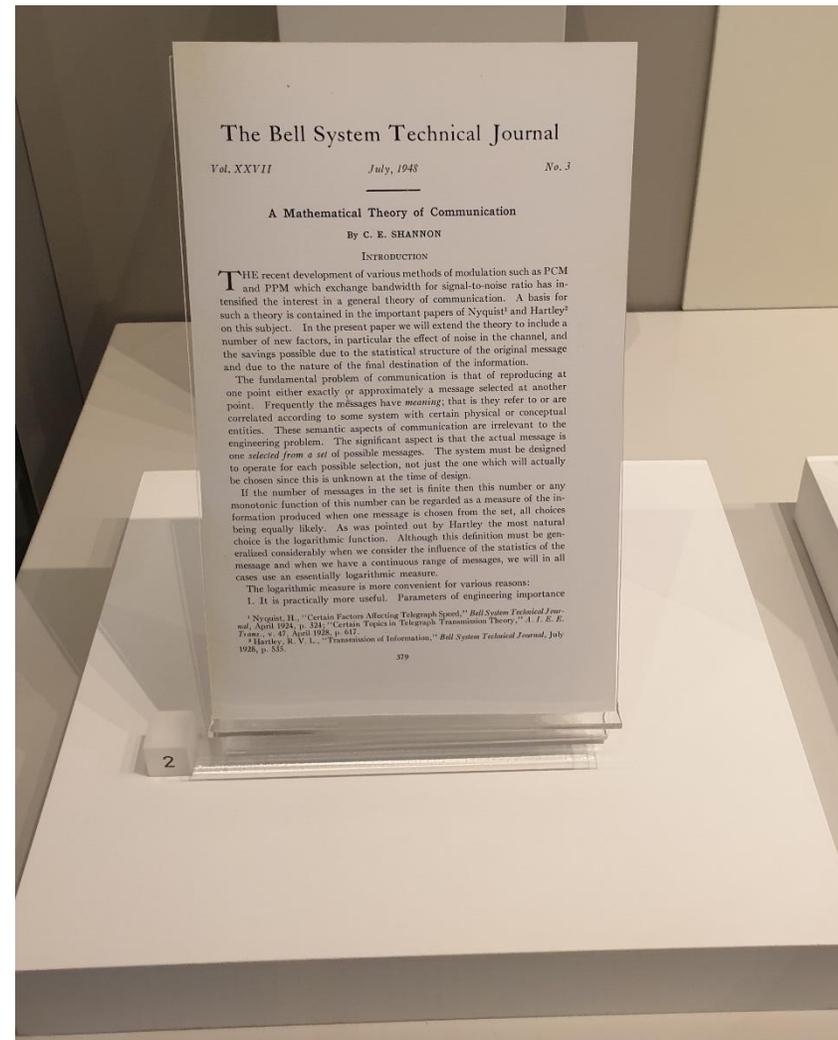
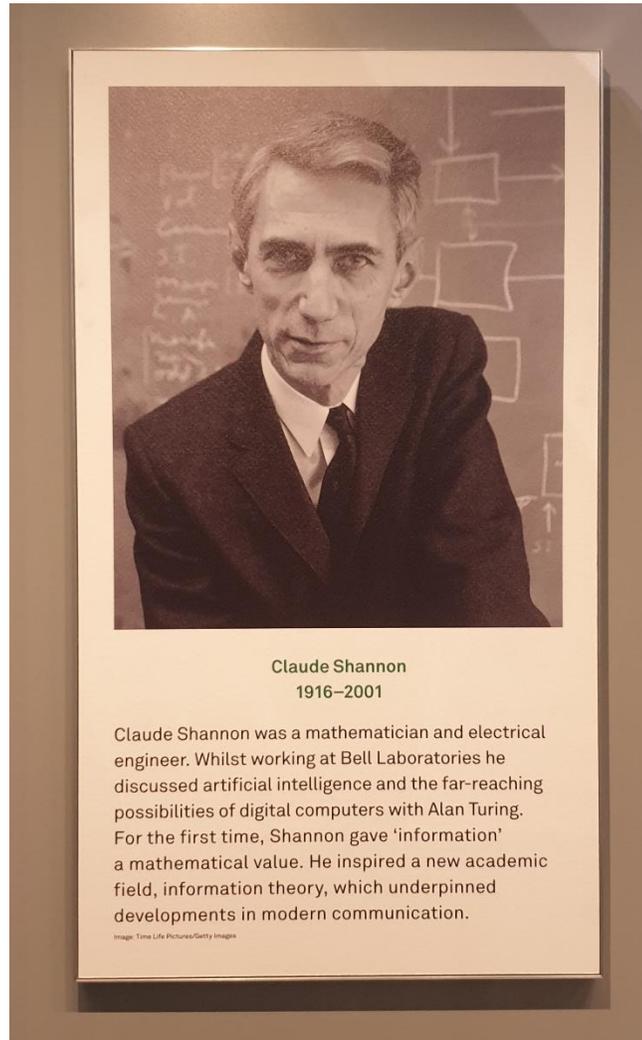
## Le principe de l'entropie maximale

- Maximiser l'entropie : une méthode générale pour sélectionner un prior en tenant compte de :
  - l'indifférence à l'égard des états de connaissance égale
  - d'informations préalables pertinentes
- Quelle forme doit avoir  $H[p]$  pour une source d'informations produisant  $N$  « mots » finis avec des probabilités respectives  $p_n$  ?
- Désidératas :
  - Si tous les mots sont équiprobables ( $\forall n, p_n = \frac{1}{N}$ ), alors  $H[p]$  doit croître avec  $N$ .
  - Si les mots sont générés en deux étapes (1 : choisir un sous-ensemble de mots ; 2 : choisir un mot dans ce sous-ensemble), alors l'entropie doit être la somme des entropies associées à chaque étape.

➔ Théorème (Shannon):  $H[p] \propto - \sum_n p_n \log_2 p_n$

# La naissance de la théorie de l'information

- Photos prises au *Science Museum* (South Kensington, Londres) en 2021 :



Pourquoi ne pas appeler ça de l'entropie ? Premièrement, un développement mathématique très semblable au vôtre existe déjà dans la mécanique statistique de Boltzmann, et deuxièmement, personne ne comprend très bien l'entropie, donc dans n'importe quelle discussion, vous serez dans une position avantageuse.

von Neumann à Shannon, à propos d'un nom pour « l'information manquante »

## Le problème du dé pipé

- Pour un dé non pipé,  $p_n = \frac{1}{6} \quad \forall n \in [1, 6]$  : le principe d'indifférence suffit.
- Maintenant supposons que la valeur moyenne après de nombreux lancers n'est pas 3,5 mais 4. Quelle est la loi de probabilité dans ce cas ?
- Il nous faut maximiser  $H[p]$  avec deux contraintes :

$$\langle n \rangle_p = \sum_{n=1}^6 n p_n = 4 \quad \text{et} \quad \sum_{n=1}^6 p_n = 1$$

- Méthode (1) : par force brute !
  - Obtenir  $p_5$  et  $p_6$  en fonction de  $p_1, p_2, p_3, p_4$ .
  - Exprimer  $H[p] = \sum_{n=1}^6 p_n \ln p_n$  en fonction de  $p_1, p_2, p_3, p_4$ .
  - Dériver et résoudre  $\frac{\partial H}{\partial p_n} = 0$  pour tout  $n \in [1, 4]$ .

## Le problème du dé pipé

- Méthode (2) : une solution plus élégante qui respecte la symétrie des variables : la méthode des **multiplicateurs de Lagrange**.

- On écrit le Lagrangien :  $\mathcal{L}[\{p_n\}, \lambda, \mu] = - \sum_{n=1}^6 p_n \ln p_n - \lambda \left( \sum_{n=1}^6 n p_n - 4 \right) - \mu \left( \sum_{n=1}^6 p_n - 1 \right)$

- Les deux contraintes sont :  $\frac{\partial \mathcal{L}}{\partial \lambda} = 0$  et  $\frac{\partial \mathcal{L}}{\partial \mu} = 0$

$$\frac{\partial \mathcal{L}}{\partial p_n} = 0 \quad \text{donne} \quad -1 - \ln p_n - \lambda n - \mu = 0$$
$$p_n = \frac{e^{-\lambda n}}{Z} \quad \text{avec} \quad \ln Z \equiv 1 + \mu$$

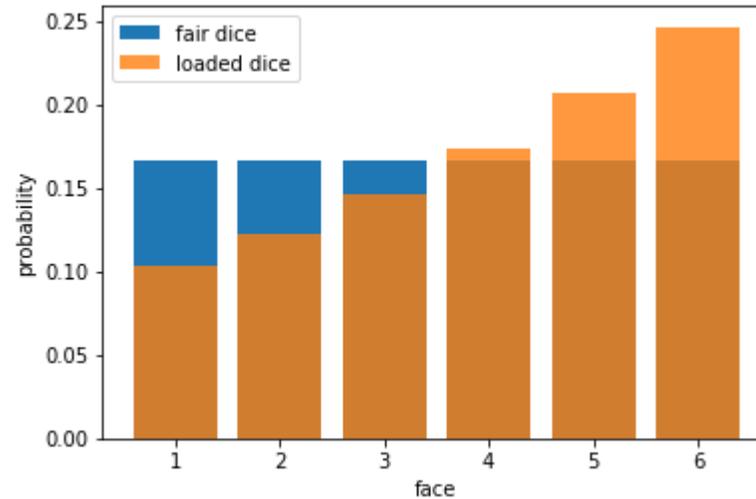
- La contrainte de normalisation impose :  $Z = \sum_{n=1}^6 e^{-\lambda n} = \frac{1 - e^{-6\lambda}}{e^\lambda - 1}$

- La contrainte sur la moyenne est obtenue en remarquant que :

$$-\frac{d \ln Z}{d\lambda} = -\frac{1}{Z} \frac{dZ}{d\lambda} = \sum_{n=1}^6 n \frac{e^{-\lambda n}}{Z} = \sum_{n=1}^6 n p_n = 4$$

- Cela donne une équation à résoudre pour  $e^\lambda$  :  $e^\lambda / (e^\lambda - 1) - 6 / (e^{6\lambda} - 1) = 4$

## Le problème du dé pipé



- Ceci est un exemple de **théorie des probabilités au-delà des statistiques bayésiennes** : nous avons obtenu numériquement les valeurs d'une distribution de probabilité, conditionnée à certaines observations, sans utiliser le théorème de Bayes.
- Analogie thermodynamique :
  - Dé non pipé = **ensemble micro-canonique**  $p_n = \frac{1}{N}$
  - Dé pipé (avec moyenne connue) = **ensemble canonique**

$$p_n = \frac{e^{-\beta E_n}}{Z} \quad \beta \equiv \frac{1}{k_B T} \quad \begin{array}{l} E_n = \text{énergie des différents états} \\ Z = \text{fonction de partition} = \text{évidence en statistiques bayésiennes} \end{array}$$

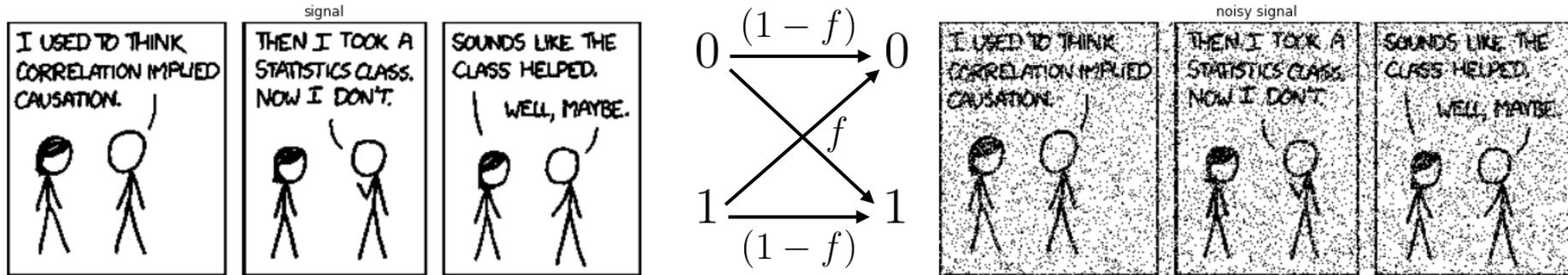
# Théorie de l'information

Notebook : <https://github.com/florent->

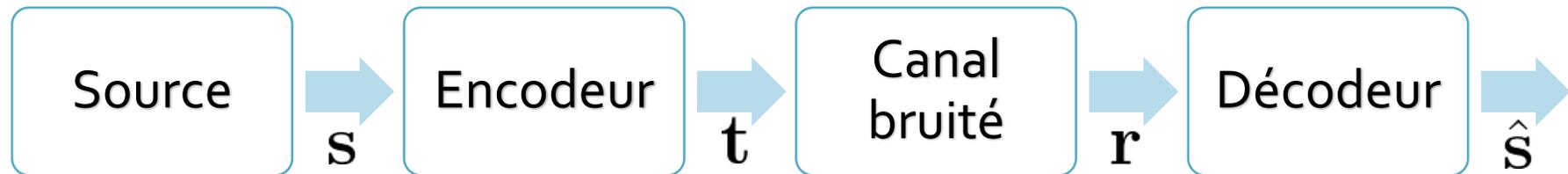
[leclercq/Bayes\\_InfoTheory/blob/master/IT\\_noisy\\_binary\\_channel.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/IT_noisy_binary_channel.ipynb)



# Le canal symétrique binaire bruité



<https://xkcd.com/552/>



Taux de transfert d'informations :  $R = \frac{\#s}{\#t} = \frac{K}{N}$

<b>s</b>	0	0	1	0	1	1	0
<b>t</b>	$\underbrace{000}$	$\underbrace{000}$	$\underbrace{111}$	$\underbrace{000}$	$\underbrace{111}$	$\underbrace{111}$	$\underbrace{000}$
<b>n</b>	000	001	000	000	101	000	000
<b>r</b>	$\underbrace{000}$	$\underbrace{001}$	$\underbrace{111}$	$\underbrace{000}$	$\underbrace{010}$	$\underbrace{111}$	$\underbrace{000}$
<b><math>\hat{s}</math></b>	0	0	1	0	0	1	0

corrected errors  
 undetected errors

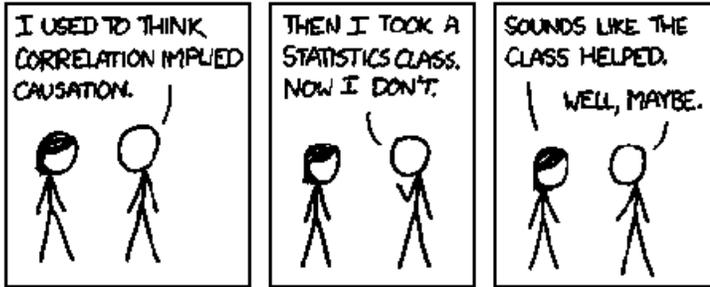
★

★

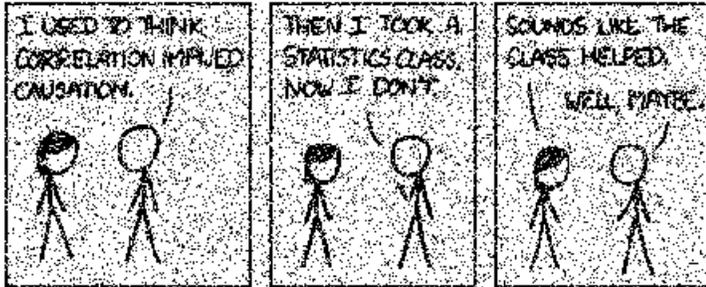
Taux de transfert d'informations :  $R[R_3] = \frac{1}{3}$

# Le code R3 : exemple

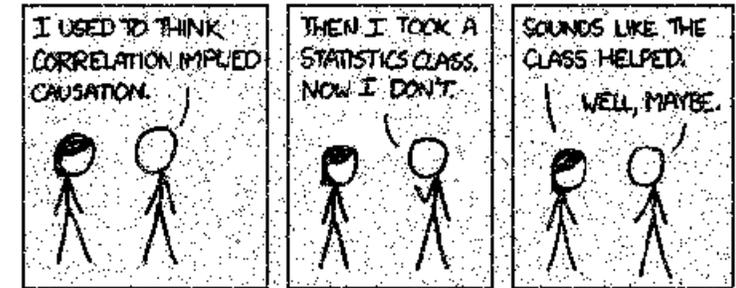
transmitted (1st)



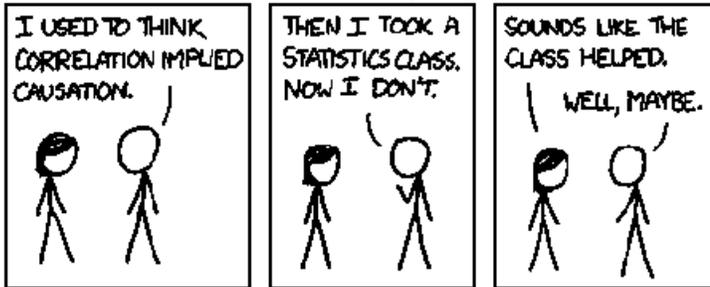
received (1st)



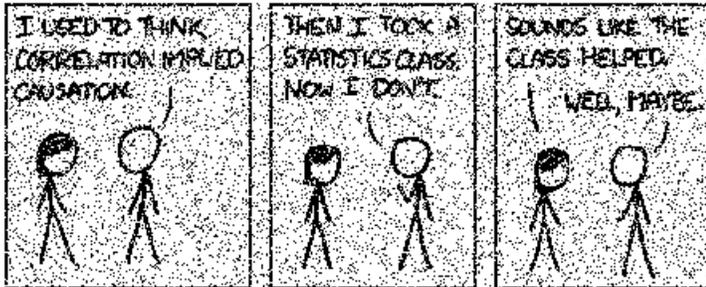
decoded



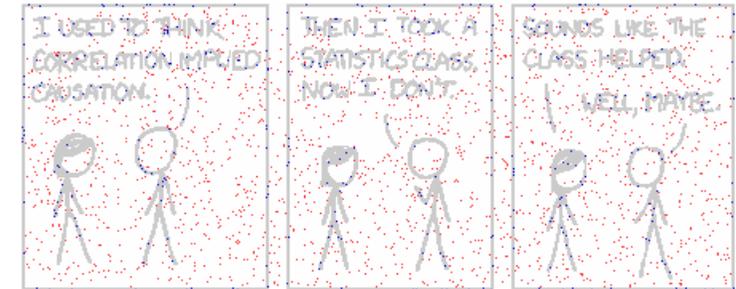
transmitted (2nd)



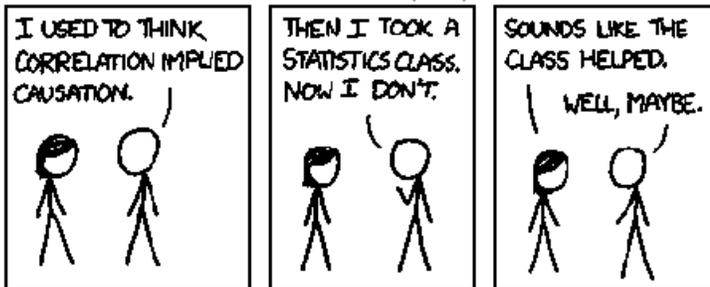
received (2nd)



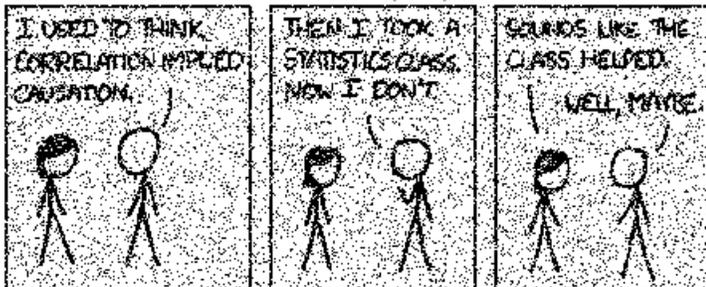
uncorrected errors



transmitted (3rd)

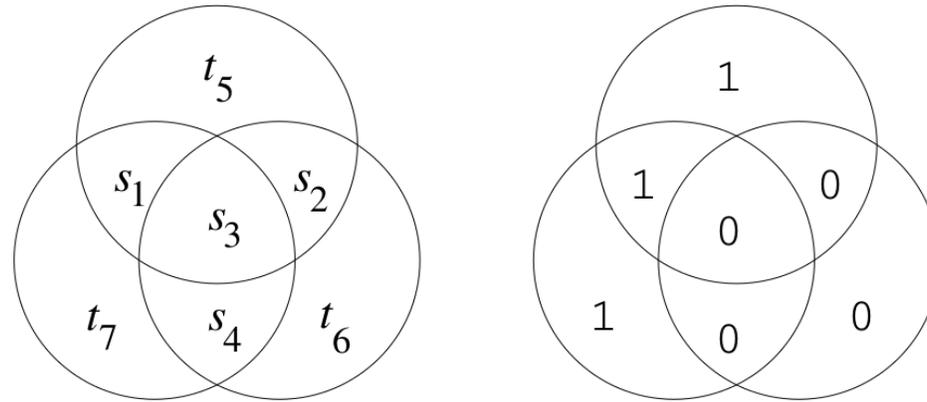


received (3rd)



## Le code de Hamming (7,4) : encodeur

- On introduit le concept de **contrôle de parité** :

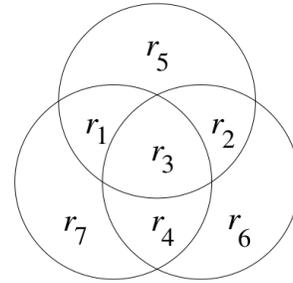


$$\mathbf{t} = \mathbf{G}\mathbf{s} \quad \mathbf{G} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

Taux de transfert d'informations :  $R[H(7,4)] = \frac{4}{7}$

# Le code de Hamming (7,4) : décodeur

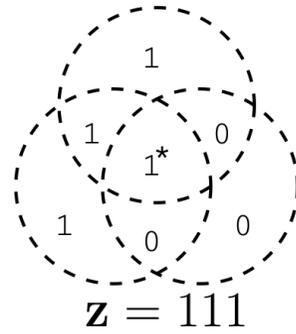
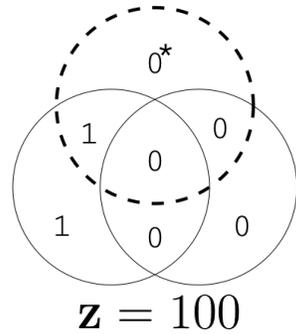
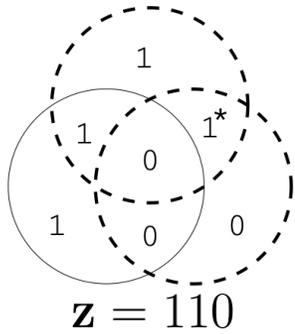
- On introduit le concept de **syndrome** :
- Algorithme pour la correction d'erreurs :



syndrome  $\rightarrow z = Hr$

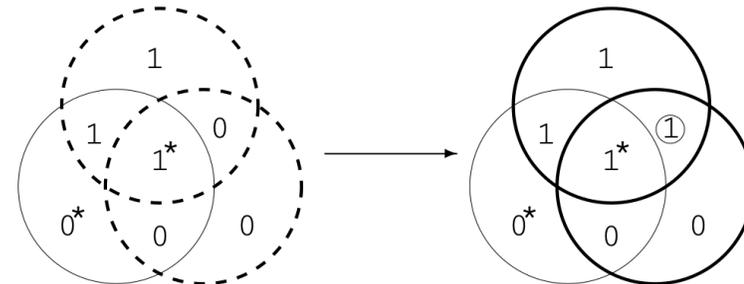
$$H = \left[ \begin{array}{cccc|ccc} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{array} \right]$$

contrôles de parité    identité



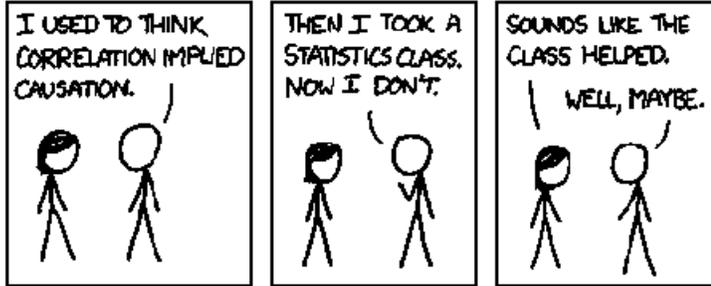
Syndrome $z$	000	001	010	011	100	101	110	111
Unflip this bit	None	$r_7$	$r_6$	$r_4$	$r_5$	$r_1$	$r_2$	$r_3$

- Exemple de correction d'erreur infructueuse :

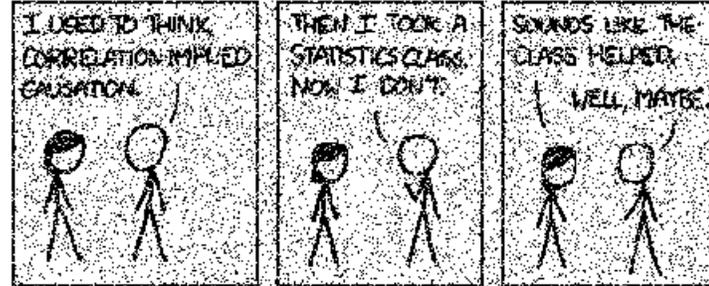


# Le code de Hamming (7,4) : exemple

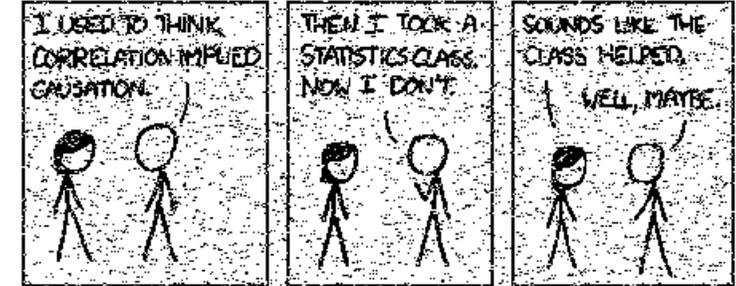
transmitted image



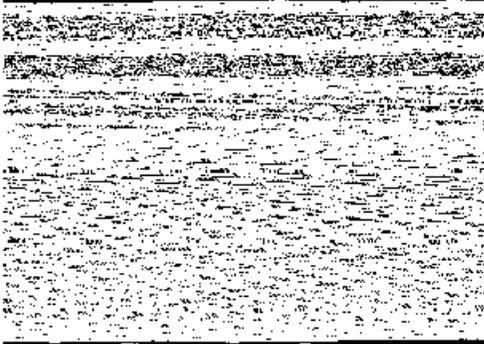
received noisy image



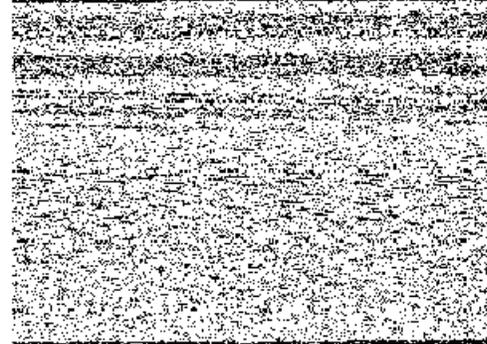
decoded



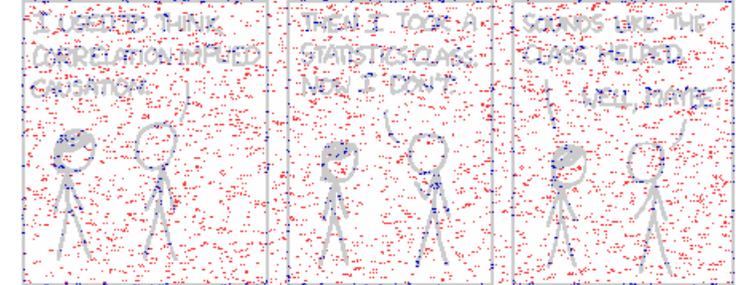
transmitted parity bits



received parity bits

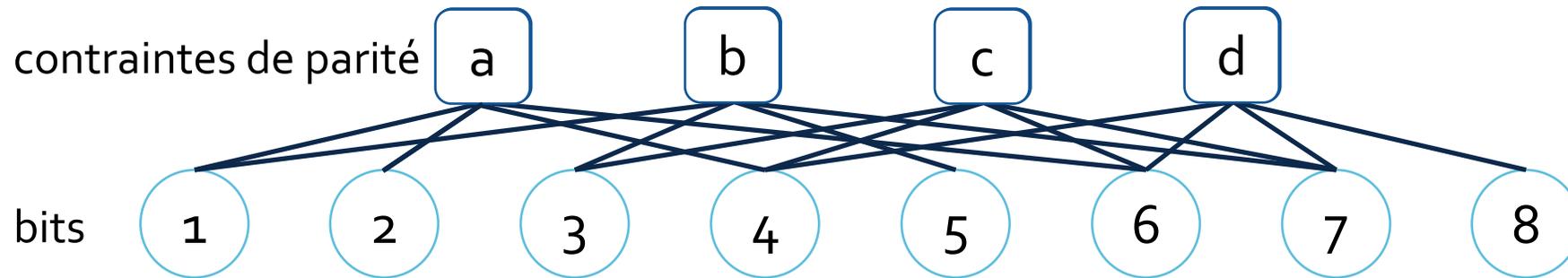


uncorrected errors



## Les codes de parité à faible densité (low-density parity-check codes – LDPC)

- Graphe de Tanner :



- $N$  bits,  $M=(1-R)N$  contraintes de parité à contrôler
  - $2^{RN}$  « mots » possibles le **dictionnaire**
  - une matrice de contrôle de parité **creuse** (*sparse matrix*)
- Décodage des codes LDPC : la théorie générale emprunte à la **physique statistique** : les spins d'Ising en interaction et l'approximation du champ moyen BP (Bethe-Peierls – Belief Propagation).

## Le théorème de Shannon pour le canal symétrique binaire bruité

- Taux de transfert d'informations  $R[R_N]$  et probabilité d'erreur  $p_b$  pour les codes à répétition (pour  $N$  impair) :

$$R[R_N] = \frac{1}{N} \quad p_b = \sum_{(N+1)/2}^N \binom{N}{n} f^n (1-f)^{N-n}$$

MacKay 2003, equation (1.24)

- Définitions :

- L'entropie (en base 2) :

$$H_2(x) \equiv x \log_2 \frac{1}{x} + (1-x) \log_2 \frac{1}{1-x}$$

- La capacité du canal bruité :

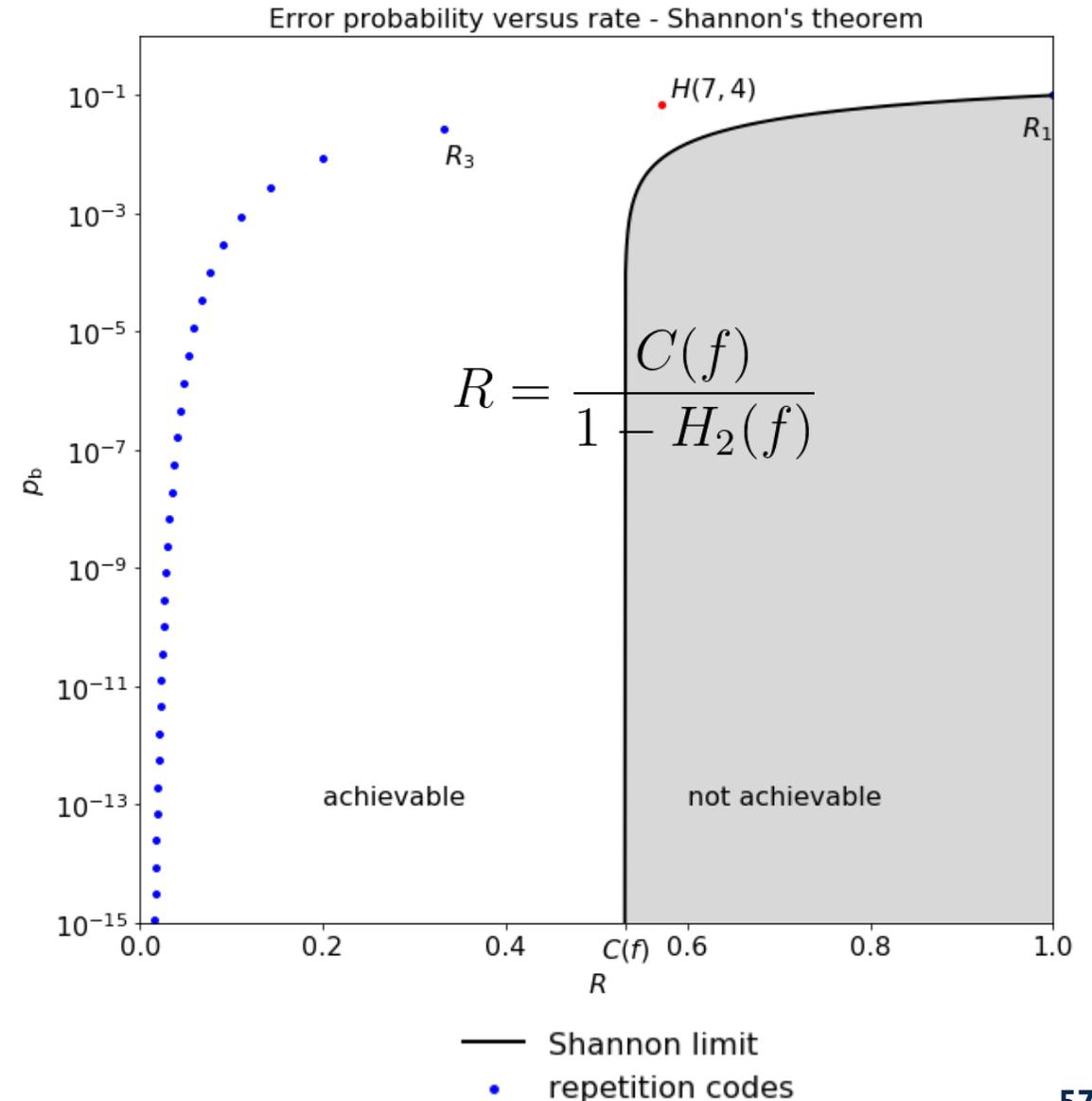
$$C(f) \equiv 1 - H_2(f)$$

- La limite de Shannon :

$$R = \frac{C(f)}{1 - H_2(f)}$$

## Le théorème de Shannon pour le canal symétrique binaire bruité

- [Théorème de Shannon](#) (1948) : les taux de transfert d'informations supérieurs à la limite de Shannon sont réalisables, les taux inférieurs à la limite de Shannon ne sont pas réalisables.
- La frontière entre régions réalisables et non-réalisables rencontre l'axe des  $R$  à une valeur non-nulle  $R = C(f)$ . On n'a pas nécessairement  $R \rightarrow 0$  si  $p_b \rightarrow 0$ .
- Les taux de transfert d'informations jusqu'à la capacité du canal  $R = C(f)$  sont réalisables avec une probabilité d'erreur  $p_b$  arbitrairement faible.



# Théorie bayésienne de la décision

Notebook : <https://github.com/florent->

[leclercq/Bayes\\_InfoTheory/blob/master/DecisionTheory.ipynb](https://github.com/florent-leclercq/Bayes_InfoTheory/blob/master/DecisionTheory.ipynb)

- La théorie bayésienne de la décision est un cadre pour la **prise de décision optimale**, étant donné un ensemble d'actions possibles et un état de connaissance incertain, représenté par une pdf  $p(\theta|I)$  (habituellement le posterior d'une inférence bayésienne préalable à la prise de décision).
- Notations :
  - $\{\theta\}$  = ensemble de paramètres (variables observées)
  - $\{a\}$  = ensemble d'actions possibles
- Hypothèse de l'utilité attendue : étant donné un ensemble de fonctions de gain  $G(a|\theta)$ , la règle de décision optimale consiste à effectuer l'action qui maximise l'utilité attendue  $U(a|I)$ , définie par
$$U(a|I) \equiv \langle G(a|\theta) \rangle_{p(\theta|I)} = \int G(a|\theta)p(\theta|I) d\theta$$
- Il faut donc effectuer l'action  $a^* = \operatorname{argmax}_a U(a|I)$ .

## Exemple : alertes bayésiennes

- On recherche un évènement  $E$ . On a accès à  $p(E|I)$  et  $p(\bar{E}|I) = 1 - p(E|I)$ .
- Il y a deux actions possibles :
  - $a_1$  = donner l'alerte
  - $a_2$  = ne rien faire

- Les fonctions d'utilité sont :

$$\begin{aligned}
 U(a_1|I) &= \overbrace{G(a_1|E)p(E|I)}^{\substack{\text{détecton correcte} \\ \text{(un « hit »)}}} + \overbrace{G(a_1|\bar{E})[1 - p(\bar{E}|I)]}^{\substack{\text{faux positif} \\ \text{(une « fausse alarme »)}}} \\
 U(a_2|I) &= \overbrace{G(a_2|E)p(E|I)}^{\substack{\text{faux négatif} \\ \text{(un « miss »)}}} + \overbrace{G(a_2|\bar{E})[1 - p(\bar{E}|I)]}^{\substack{\text{rejet correcte}}}
 \end{aligned}$$

- Un choix typique de fonctions de gain :
 
$$\begin{aligned}
 G(a_1|E) &= G - C & G(a_1|\bar{E}) &= -C \\
 G(a_2|E) &= 0 & G(a_2|\bar{E}) &= 0
 \end{aligned}$$

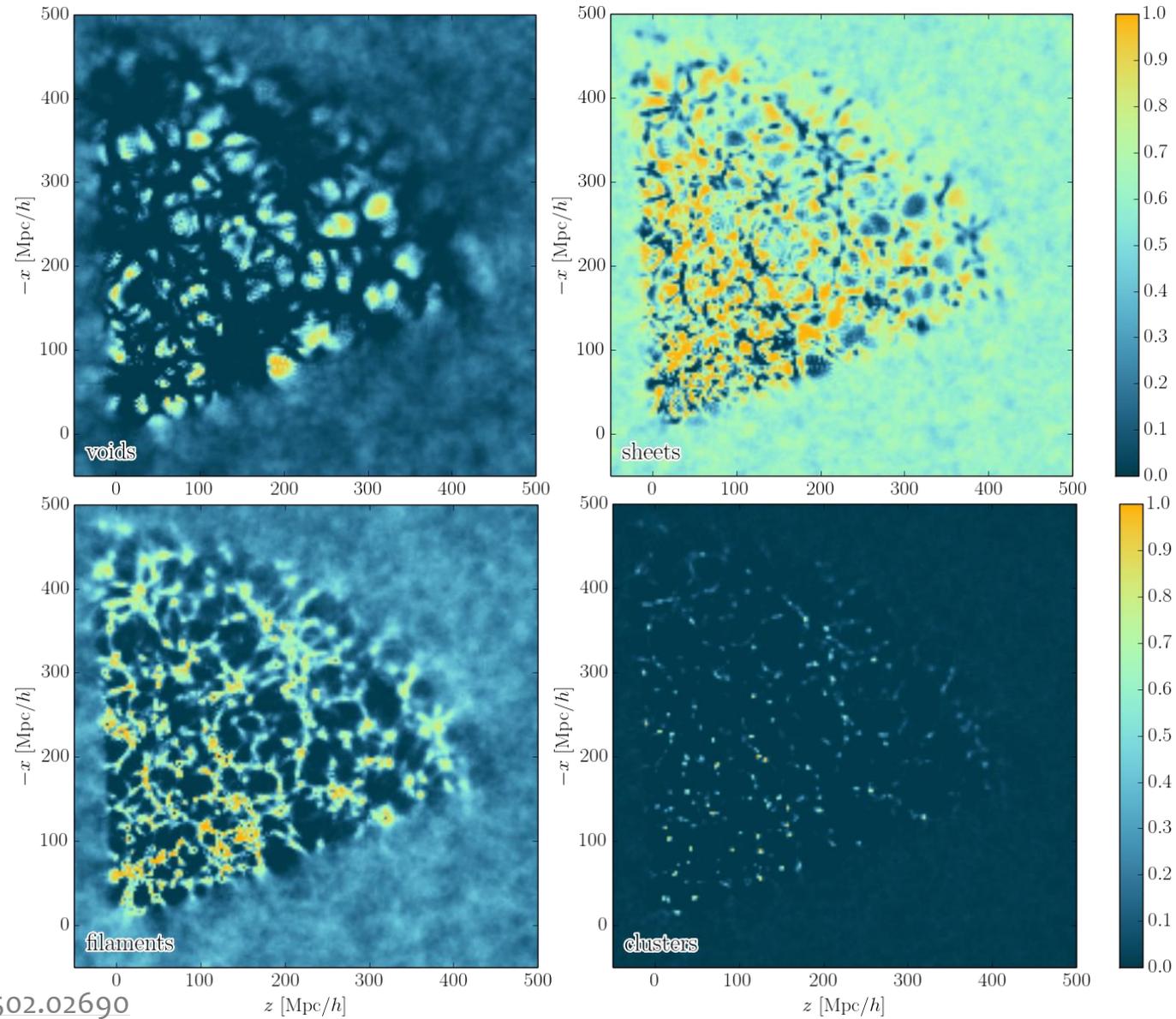
le gain attendu pour une détection
le coût de donner une alerte

- On a donc  $U(a_1|I) = p(E|I)(G - C) + [1 - p(E|I)](-C)$

$$U(a_2|I) = 0$$

➔ il faut donner l'alerte si et seulement si  $p(E|I) \geq \frac{C}{G}$

# Classification des structures dans la toile cosmique



FL, Jasche & Wandelt 2015a, 1502.02690

## Une règle décisionnelle pour la classification des structures dans la toile cosmique

- Une espace de « paramètres d'entrée » :

$$\{T_0 = \text{void}, T_1 = \text{sheet}, T_2 = \text{filament}, T_3 = \text{cluster}\}$$

- Un espace de « décisions possibles » :

$$\{a_0 = \text{“decide void”}, a_1 = \text{“decide sheet”}, a_2 = \text{“decide filament”}, \\ a_3 = \text{“decide cluster”}, a_{-1} = \text{“do not decide”}\}$$

- C'est donc un problème de [théorie bayésienne de la décision](#) : il faut effectuer l'action qui maximise l'utilité attendue

$$U(a_j(\vec{x}_k)|d) = \sum_{i=0}^3 G(a_j|T_i) \mathcal{P}(T_i(\vec{x}_k)|d)$$

- Comment choisir les fonctions de gain ?

## Un jeu de hasard avec l'Univers

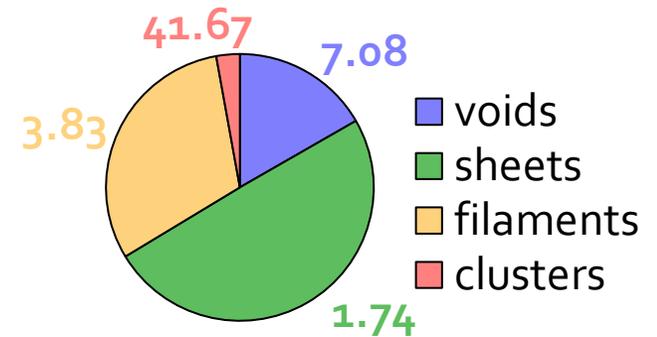
- Une proposition :

$$G(a_j | \mathbf{T}_i) = \begin{cases} \frac{1}{\mathcal{P}(\mathbf{T}_i)} - \alpha & \text{if } j \in \llbracket 0, 3 \rrbracket \text{ and } i = j & \text{« gagner »} \\ -\alpha & \text{if } j \in \llbracket 0, 3 \rrbracket \text{ and } i \neq j & \text{« perdre »} \\ 0 & \text{if } j = -1. & \text{« ne pas jouer »} \end{cases}$$

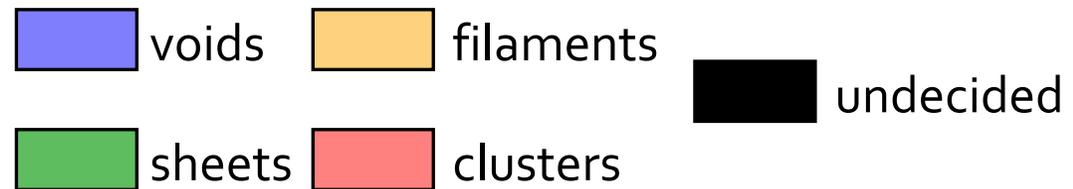
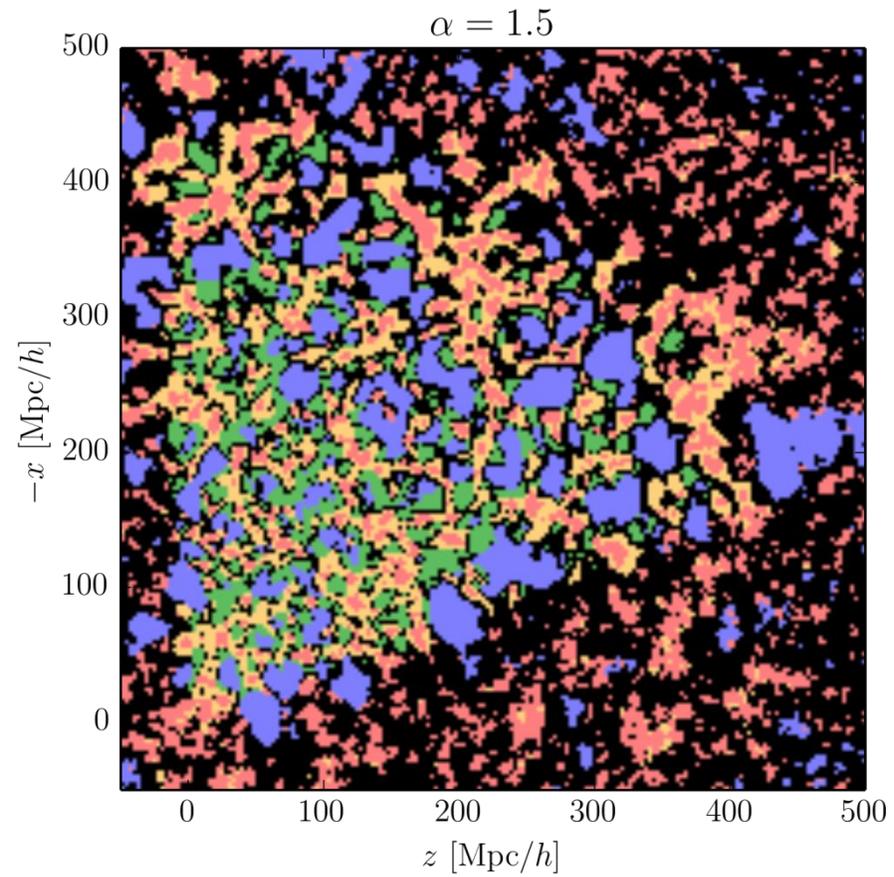
- En l'absence de données, l'utilité attendue est :

$$U(a_j) = 1 - \alpha \quad \text{if } j \neq -1 \quad \text{« jouer au jeu »}$$
$$U(a_{-1}) = 0 \quad \text{« ne pas jouer au jeu »}$$

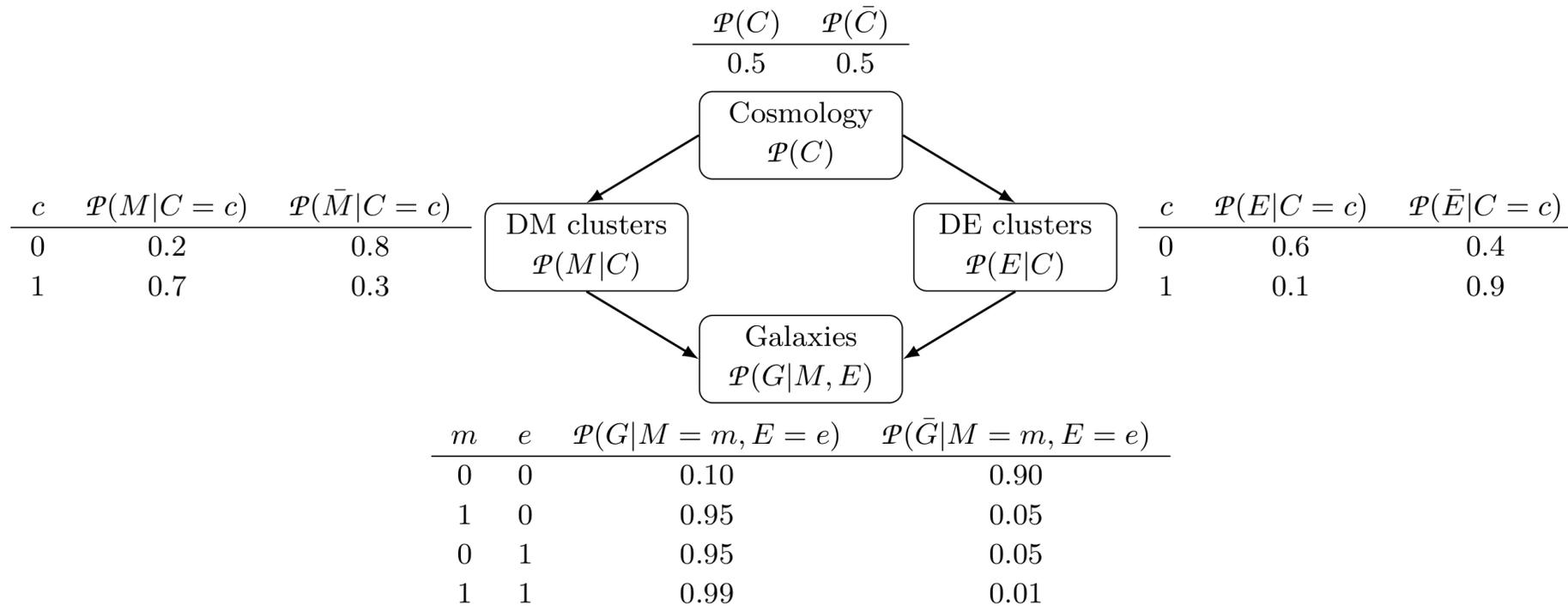
- Pour  $\alpha = 1$ , c'est un **jeu équilibré** (*fair game*)  $\Rightarrow$  on va donc toujours jouer  $\Rightarrow$  on construit une « **carte spéculative** » des structures cosmiques.
- Des valeurs  $\alpha > 1$  représentent une **aversion pour le risque**  $\Rightarrow$  on ne joue que si les données sont informatives  $\Rightarrow$  on construit des « **cartes conservatives** » des structures cosmiques.



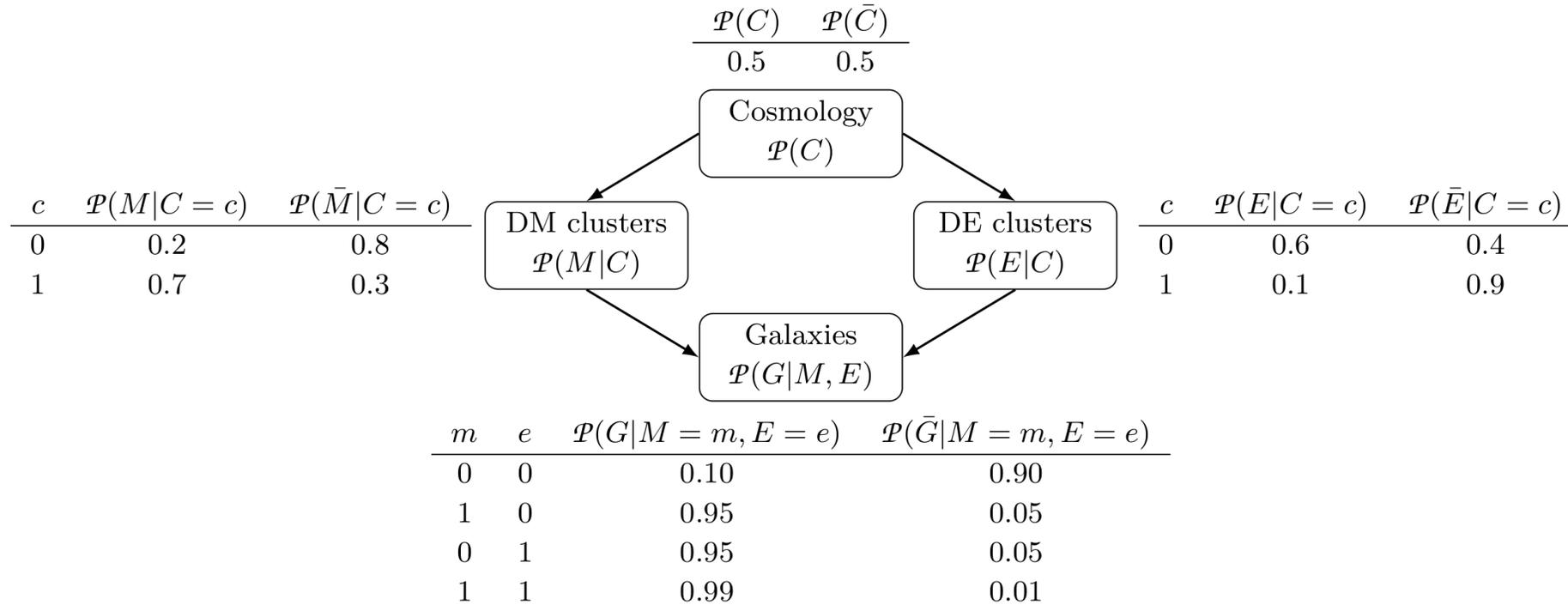
# Jouons au jeu...



# Réseaux bayésiens, modèles bayésiens hiérarchiques et Bayésianisme empirique



- Les réseaux bayésiens sont des modèles probabilistes (avec une représentation graphique), comprenant :
  - Une **graphe orienté acyclique** (directed acyclic graph – DAG)
  - À chacun des nœuds, les **distributions de probabilités conditionnelles**



- Le graphe permet de simplifier l'écriture de la probabilité jointe de toutes les variables :

$$p(C, M, E, G) = p(C) p(E|C) p(M|C, \cancel{E}) p(G|\cancel{C}, M, E)$$

$$p(C, M, E, G) = p(C) p(E|C) p(M|C) p(G|M, E)$$

- Inférence :

$$p(M|G) = \frac{p(M,G)}{p(G)} = \frac{\sum_{c,e} p(C=c, M=1, E=e, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.4313}{0.70305} \approx 0.6135$$

$$p(E|G) = \frac{p(E,G)}{p(G)} = \frac{\sum_{c,m} p(C=c, M=m, E=1, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.3363}{0.70305} \approx 0.4783$$

$$p(\bar{M}, \bar{E}|G) = \frac{p(\bar{M}, \bar{E}, G)}{p(G)} = \frac{\sum_c p(C=c, M=0, E=0, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.0295}{0.70305} \approx 0.0420$$

- Prédiction :

$$p(G|C) = \frac{p(G,C)}{p(C)} = \frac{\sum_{m,e} p(C=1, M=m, E=e, G=1)}{p(C=1)} = 0.7233$$

$$p(E|M, G) = \frac{p(E, M, G)}{p(M, G)} = \frac{\sum_c p(C=c, M=1, E=1, G=1)}{\sum_{c,e} p(C=c, M=1, E=e, G=1)} = \frac{0.09405}{0.4313} \approx 0.2181$$

$$p(E|G) = \frac{p(E, G)}{p(G)} = \frac{\sum_{c,m} p(C=c, M=m, E=1, G=1)}{\sum_{c,m,e} p(C=c, M=m, E=e, G=1)} = \frac{0.3363}{0.70305} \approx 0.4783$$

- On a donc à la fois :

$$p(E|M) = p(E)$$

$$p(E|M, G) < p(E|G)$$

- $G$  est appelé un **collisionneur**. Ce phénomène est appelé « **biais de collision** » ou « explaining away » : deux causes rentrent en compétition pour expliquer le même effet.
- Cas particulier : le « **biais de sélection** » ou « **paradoxe de Berkson** » :

$$0 < p(A) < 1; \quad 0 < p(B) < 1; \quad p(A|B) = p(A)$$

$$C = A + B \quad \rightarrow \quad \begin{aligned} p(A|B, C) &< p(A|C) \\ p(A|\bar{B}, C) &= 1 > p(A|C) \end{aligned}$$

## Le biais de Malmquist

- [Le biais de Malmquist](#) (1925) : dans les relevés astronomiques limités par la magnitude des objets, les objets éloignés sont préférentiellement détectés s'ils sont intrinsèquement brillants.



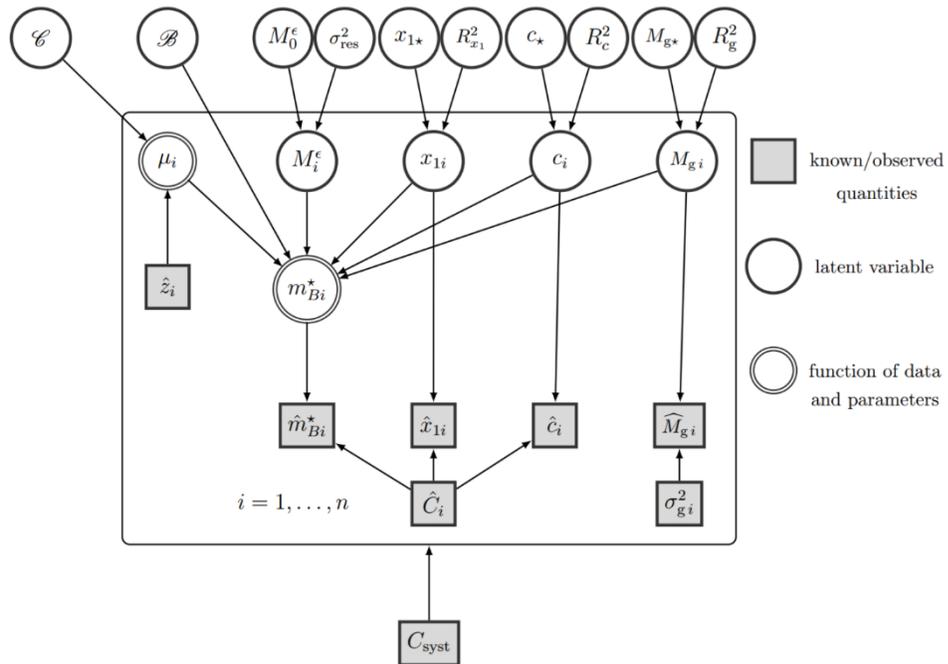
Gunnar Malmquist  
(1893-1982)



$$0 < p(A) < 1; \quad 0 < p(B) < 1; \quad p(A|B) = p(A)$$
$$C = A + B \quad \Rightarrow \quad p(A|\bar{B}, C) = 1 > p(A|C)$$

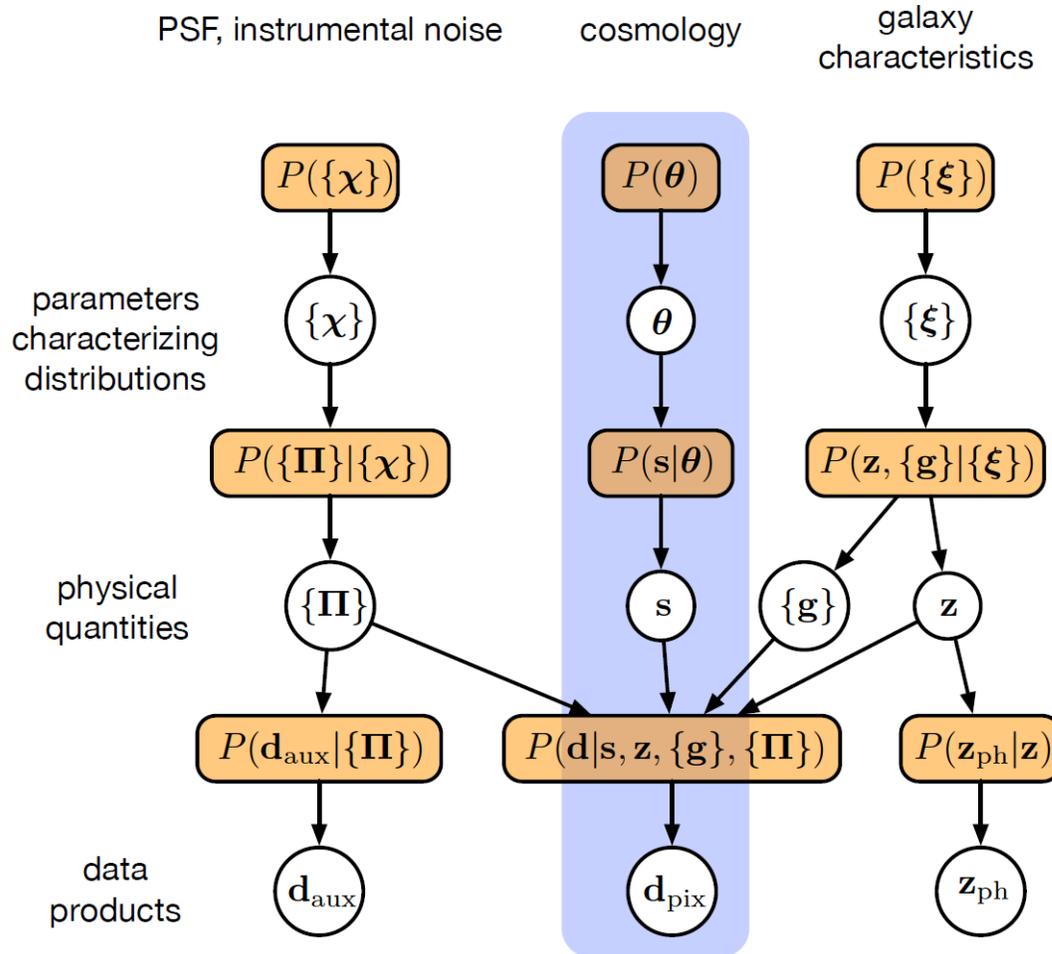
déTECTÉ      lumineux      proche

# Exemple de réseau bayésien : supernovæ (BAHAMAS)



Parameter	Notation and Prior Distribution
Cosmological parameters	
Matter density parameter	$\Omega_m \sim \text{UNIFORM}(0, 2)$
Cosmological constant density parameter	$\Omega_\Lambda \sim \text{UNIFORM}(0, 2)$
Dark energy EOS	$w \sim \text{UNIFORM}(-2, 0)$
Hubble parameter	$H_0/\text{km/s/Mpc} = 67.3$
Covariates	
Coefficient of stretch covariate	$\alpha \sim \text{UNIFORM}(0, 1)$
Coefficient of color covariate	$\beta$ (or $\beta_0$ ) $\sim \text{UNIFORM}(0, 4)$
Coefficient of interaction of color correction and $z$	$\beta_1 \sim \text{UNIFORM}(-4, 4)$
Jump in coefficient of color covariate	$\Delta\beta \sim \text{UNIFORM}(-1.5, 1.5)$
Redshift of jump in color covariate	$z_t \sim \text{UNIFORM}(0.2, 1)$
Coefficient of host galaxy mass covariate	$\gamma \sim \text{UNIFORM}(-4, 4)$
Population-level distributions	
Mean of absolute magnitude	$M_0^\epsilon \sim \mathcal{N}(-19.3, 2^2)$
Residual scatter after corrections	$\sigma_{\text{res}}^2 \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of absolute magnitude, low galaxy mass	$M_0^{\text{lo}} \sim \mathcal{N}(-19.3, 2^2)$
SD of absolute magnitude, low galaxy mass	$\sigma_{\text{res}}^{\text{lo}^2} \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of absolute magnitude, high galaxy mass	$M_0^{\text{hi}} \sim \mathcal{N}(-19.3, 2^2)$
SD of absolute magnitude, high galaxy mass	$\sigma_{\text{res}}^{\text{hi}^2} \sim \text{INV GAMMA}(0.003, 0.003)$
Mean of stretch	$x_{1*} \sim \mathcal{N}(0, 10^2)$
SD of stretch	$R_{x_1} \sim \text{LOG UNIFORM}(-5, 2)$
Mean of color	$c_* \sim \mathcal{N}(0, 1^2)$
SD of color	$R_c \sim \text{LOG UNIFORM}(-5, 2)$
Mean of host galaxy mass	$M_{g*} \sim \mathcal{N}(10, 100^2)$
SD of host galaxy mass	$R_g \sim \text{LOG UNIFORM}(-5, 2)$

# Exemple de réseau bayésien : lentillage gravitationnel faible



Can include:

- Mask
- Intrinsic alignments
- Baryon feedback
- Shape measurement
- Photometric redshifts

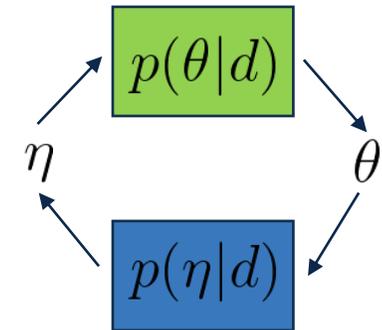
Le bayésianisme empirique :  
une alternative au principe d'entropie maximale pour choisir un prior

$$p(\theta|d) \propto p(d|\theta) \overset{\text{prior}}{p(\theta|\eta)} \overset{\text{hyperprior}}{p(\eta)}$$

$$\underline{p(\theta|d)} = \int p(\theta|\eta, d) p(\eta|d) d\eta = \int \frac{p(d|\theta) p(\theta|\eta)}{p(d|\eta)} \underline{p(\eta|d)} d\eta$$

$$\underline{p(\eta|d)} = \int p(\eta|\theta) \underline{p(\theta|d)} d\theta$$

➔ Schéma itératif (analogue à un échantillonneur de Gibbs)



- Le **bayésianisme empirique** (Empirical Bayes) est une troncature de ce schéma après quelques étapes (souvent une seule).

- Cas particulier :  $p(\eta|d) \approx \delta_D(\eta - \eta^*(d)) \Rightarrow \underline{p(\theta|d)} \approx \frac{p(d|\theta) p(\theta|\eta^*)}{p(d|\eta^*)}$

Cet algorithme s'appelle l'**algorithme Espérance-Maximisation** (Expectation-Maximisation – EM) en machine learning, data mining.