

# Chapter 4: Statistical uncertainty

UE Computational statistics

Pierre Pudlo

Aix-Marseille Université / Faculté des Sciences

# About this chapter

## Objectives:

- How the computer can be harnessed to obtain measures of statistical uncertainty such as standard errors, confidence intervals, etc.

## Content:

- Fisher and the maximum likelihood estimator (MLE)
- The bootstrap
- MCMC methods to approximate the posterior distribution in Bayesian statistics



### Warning

This chapter is **much longer** than the first chapters.

# 1. Introduction

## Statistical uncertainty

Assess the accuracy or the uncertainty of something computed from data: estimate of a parameter, predictions on new individuals, etc.

In linear regression, e.g., the response of individual  $i$  is

$$Y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \varepsilon_i$$

where  $\varepsilon_i$  is the error term, independent of the covariates  $X_{ij}$ 's, with  $\mathbb{E}(\varepsilon_i) = 0$  and  $\mathbb{V}(\varepsilon_i) = \sigma^2$ .

- Estimation uncertainty: true  $\beta_j$ 's  $\neq$  estimates,...
- Prediction uncertainty: exact prediction impossible because  $\sigma^2 > 0$

# 1.1 Statistics (revisited): the statistical model

- All observed numbers are modeled by their own random variables (rv's)
- In linear regression, e.g., we introduce many rv's:

$$D = \begin{pmatrix} Y_1 & X_{11} & \cdots & X_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ Y_n & X_{n1} & \cdots & X_{np} \end{pmatrix}$$

- The individuals  $D_i = (Y_i \ X_{i1} \ \cdots \ X_{ip})$  are independent and identically distributed (iid) random vectors
- We might introduce a new individual  $D_{n+1} = (Y_{n+1} \ X_{n+1,1} \ \cdots \ X_{n+1,p})$  and assume that  $Y_{n+1}$  is not observed

- In **logistic regression**, the dataset is modeled likewise, but the response  $Y_i$  is binary, i.e.,  $Y_i \in \{0, 1\}$  and we assume that

$$[Y_i|X_i] \sim \text{Bernoulli}\left(p(X_i)\right), \quad \text{with } p(X_i) = \frac{1}{1 + \exp\left(-\beta_0 - \sum_j \beta_j X_{ij}\right)}.$$

- In a **Poisson regression**, the response  $Y_i$  is a count, i.e.,  $Y_i \in \{0, 1, 2, \dots\}$  and we assume that

$$[Y_i|X_i] \sim \text{Poisson}\left(\lambda(X_i)\right), \quad \text{with } \lambda(X_i) = \exp\left(\beta_0 + \sum_j \beta_j X_{ij}\right).$$

- In a **Gaussian mixture model**, the model  $Y_1, \dots, Y_n$  is composed of  $n$  iid random vectors and we assume that

$$Y_i \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma_k), \quad \text{with } \sum_{k=1}^K \pi_k = 1.$$

# 1.2 Frequentist statistics

- In the frequentist paradigm, the parameters  $\theta$  of the statistical model are fixed but **unknown**
  - E.g.,  $\beta_0, \dots, \beta_p$  in regression problems,  $\sigma^2 > 0$  in linear regression,
  - $\pi_k$ 's,  $\mu_k$ 's and  $\Sigma_k$ 's in the Gaussian mixture model, etc.
- The observed dataset  $d$  is a realization of the random dataset  $D$  with distribution  $P_\theta$
- Sometimes, rv's are **not observed**, e.g.,  $Y_{n+1}$  in linear regression
- Two type of problems:
  - **Estimation**: approximate the unknown parameters  $\theta$  or a function of the unknown parameters  $\psi = g(\theta)$  (E.g.,  $\psi = \beta_0$  in linear regression, or  $\psi = \beta_0 + \beta_1 x_{n+1,1} + \dots + \beta_p x_{n+1,p}$  when  $x_{n+1,j}$ 's are known)
  - **Prediction**: approximate the unobserved value of a rv (E.g.,  $Y_{n+1}$  in linear regression) or its distribution

- The **uncertainty** arises from the fact that the observed dataset  $d$  is a single realization of the random dataset  $D$
- **Estimation procedures** are functions of the observed dataset  $d$ , and maybe a tuning parameters  $\lambda$ .
  - E.g., the LASSO estimate of  $\beta = (\beta_0, \dots, \beta_p)$  in linear regression is

$$\hat{\beta}(\lambda) = \underset{b \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \left[ \frac{1}{2n} \sum_{i=1}^n \left( y_i - b_0 - \sum_{j=1}^p b_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |b_j| \right]$$

- Such procedures can be seen as a **deterministic function** of the observed dataset  $d$  and the tuning parameter  $\lambda$ :

$$\hat{\beta}(\lambda) = s(d, \lambda)$$

where  $s$  is the function that implements the LASSO algorithm in Python, R, etc.

## Uncertainty measures of a parameter estimate

- **Standard errors:** an estimate of the standard deviation of  $s(D, \lambda)$
  - **Confidence intervals:** range of values that are likely to contain the true value of a parameter
- 
- To get formula for uncertainty measures, we need to resort to **probability calculus** to understand the distribution of  $s(D, \lambda)$  when  $D \sim P_\theta$ .
  - Moreover, we need **the distribution whatever the unknown value of the parameter  $\theta$**  is, i.e., we need to consider the distribution of  $s(D, \lambda)$  when  $D \sim P_\theta$  for all  $\theta$ .
  - Often in the frequentist paradigm,
    - the **estimator**  $s(D, \lambda)$  is denoted  $\widehat{\Theta}$  and
    - its observed value  $\widehat{\theta} = s(d, \lambda)$  is the **estimate**.

- How to **compute the distribution** of  $s(D, \lambda)$ ?
  - Very specific cases: the distribution of  $s(D, \lambda)$  can be **computed exactly** (Student, Fisher, Gaussian linear model,...)
  - Asymptotic results: the distribution of  $s(D, \lambda)$  can be **approximated** by a known distribution when  $n$  is large...
  - Monte Carlo methods: **the Bootstrap**, see next Section



Consider now that we want to estimate  $\psi = g(\theta)$ , and a statistic  $s(D, \lambda)$  that estimates  $\psi$ :

- $\psi$  is a given coordinate of  $\theta$
- or  $\psi = \beta_0 + \sum_j \beta_j x_{n+1,j}$  in (generalized) linear models
- or ...

## Confidence intervals

- A confidence interval of level  $(1 - \alpha)$  is an interval  $[a(D), b(D)]$  whose bounds are two statistics such that

$$\forall \theta, \quad \mathbb{P}_\theta \left( a(D) \leq g(\theta) \leq b(D) \right) \geq 1 - \alpha.$$

- For a large proportion of datasets  $D$  drawn from the model, the true value of the parameter is in the confidence interval
- The confidence interval is a **random interval** that depends on the dataset  $D$

- Once  $D$  is fixed to  $d$ ,  $a(d) \leq g(\theta) \leq b(d)$  is either true or false, but is no longer an event with a probability  $\in (0, 1)$
- We have to pray that it is true for the observed dataset  $d$ ! (We do not even know whether the data have been generated by the model or not)

# 1.3 Bayesian statistics

- In the Bayesian paradigm, **the parameters** of the statistical model **are also random variables**
  - E.g.,  $\beta_0, \dots, \beta_p$  in regression problems, and additionally  $\sigma^2 > 0$  in linear regression,
  - $\pi_k$ 's,  $\mu_k$ 's and  $\Sigma_k$ 's in the Gaussian mixture model, etc.
- The **randomness of the parameters** is a way of **modelling** the uncertainty about the values of the parameters, in other word **our belief about their values**
- To stress that the parameters are now rv's, we denote them  $\Theta$  instead of  $\theta$
- Likewise, the observed dataset  $d$  is a realization of the random dataset  $D$
- The **statistical model** is now the conditional distribution of the data given the parameter:  $[D|\Theta = \theta] \sim P_\theta$

- The **prior distribution** of the parameter:
  - is the marginal distribution of  $[\Theta]$ :
  - models our **belief** about the values of the parameters **before observing the data**
- The **posterior distribution** of the parameter is the conditional distribution of the parameter given the data:
  - it is the distribution of  $[\Theta | D = d]$
  - it takes into account the **actual data**  $d$  we have observed **to update our belief** about the values of the parameters
  - models our **belief** about the values of the parameters **after observing the data**

In the continuous case,

- the **prior density**:  $\pi(\theta)$
- the **likelihood**:  $f(d|\theta)$
- the **posterior density**:

$$\pi(\theta|d) = \frac{\pi(\theta)f(d|\theta)}{f_\pi(d)} \propto \pi(\theta)f(d|\theta), \quad \text{where } f_\pi(d) = \int \pi(\theta')f(d|\theta') d\theta'$$

- the **marginal likelihood**, or the **evidence**  $f_\pi(d)$ 
  - is the normalizing constant of the posterior distribution
  - is difficult to compute in general
  - is the probability of observing  $d$  under the Bayesian model (prior + likelihood)
  - can be used for Bayesian model selection, or for averaging over models

- The posterior distribution (our belief after the data) can be summarized:
  - the **posterior mean**:  $\mathbb{E}[\Theta | D = d]$
  - the **posterior variance**:  $\mathbb{V}[\Theta | D = d]$  and standard deviation
  - the **posterior median**, and other quantiles
- The posterior mean and median can serve as **estimates** of the true parameter
- The posterior variance, standard deviation and quantiles can serve as **measures of uncertainty** of the estimate
- The **main difficulties in Bayesian statistics** are:
  - **set the prior** to reflect our prior belief
  - **compute the posterior** distribution

# Basic example

## The Beta distribution $B(\alpha, \beta)$

- The Beta distribution is a continuous distribution on the interval  $(0, 1)$
- Has density function

$$f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \quad \text{where } B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

- If  $X \sim B(\alpha, \beta)$ , then

$$\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta}, \quad \mathbb{V}[X] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$

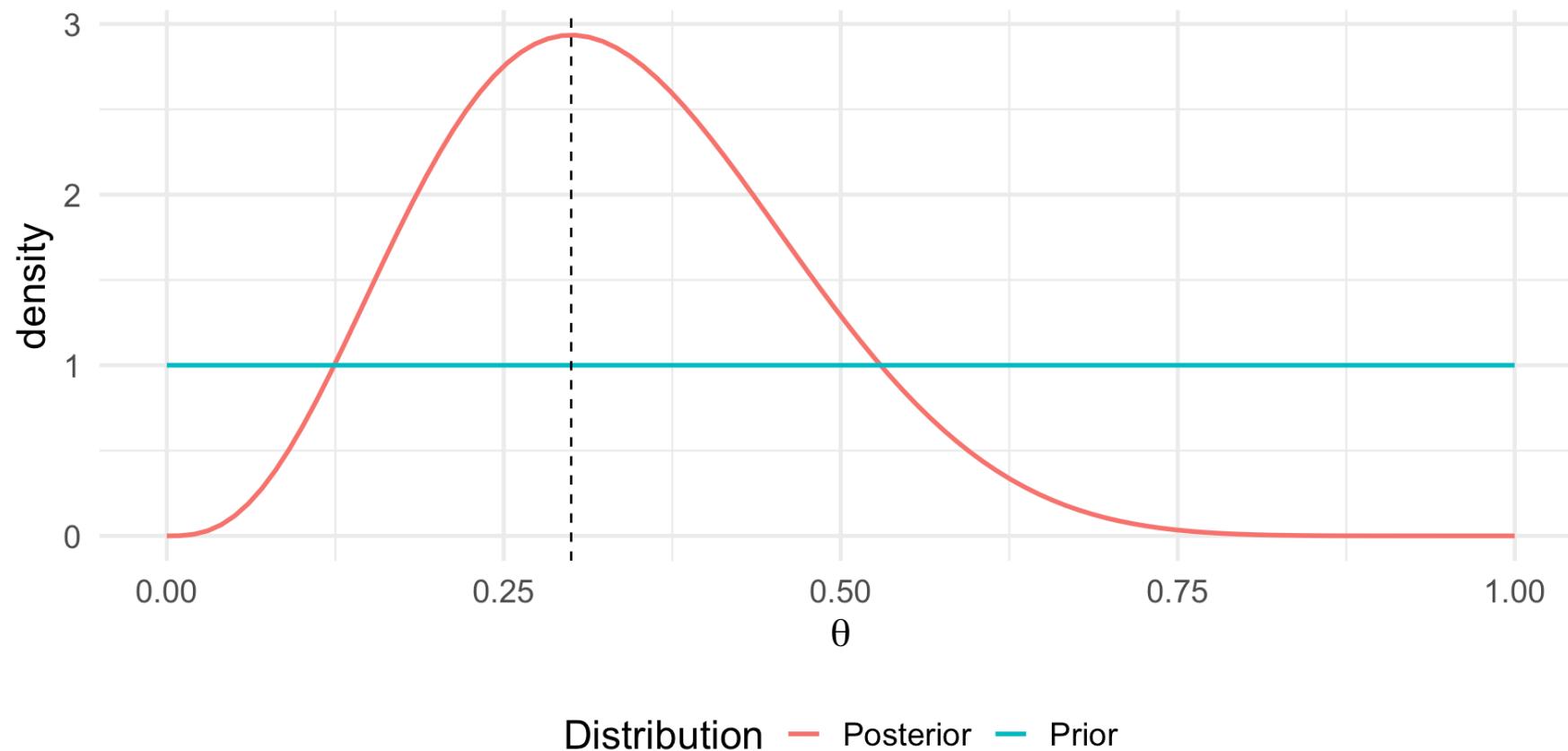
- $[D|\Theta = \theta] \sim \text{Binomial}(n, \theta)$  with unknown  $\theta \in (0, 1)$
- The likelihood of  $\theta$  is

$$f(d|\theta) = \binom{n}{d} \theta^d (1-\theta)^{n-d}$$

- The **prior** distribution of  $\Theta$  is  $B(\alpha, \beta)$
- The **posterior** distribution of  $\Theta$  is  $B(\alpha + d, \beta + n - d)$ :

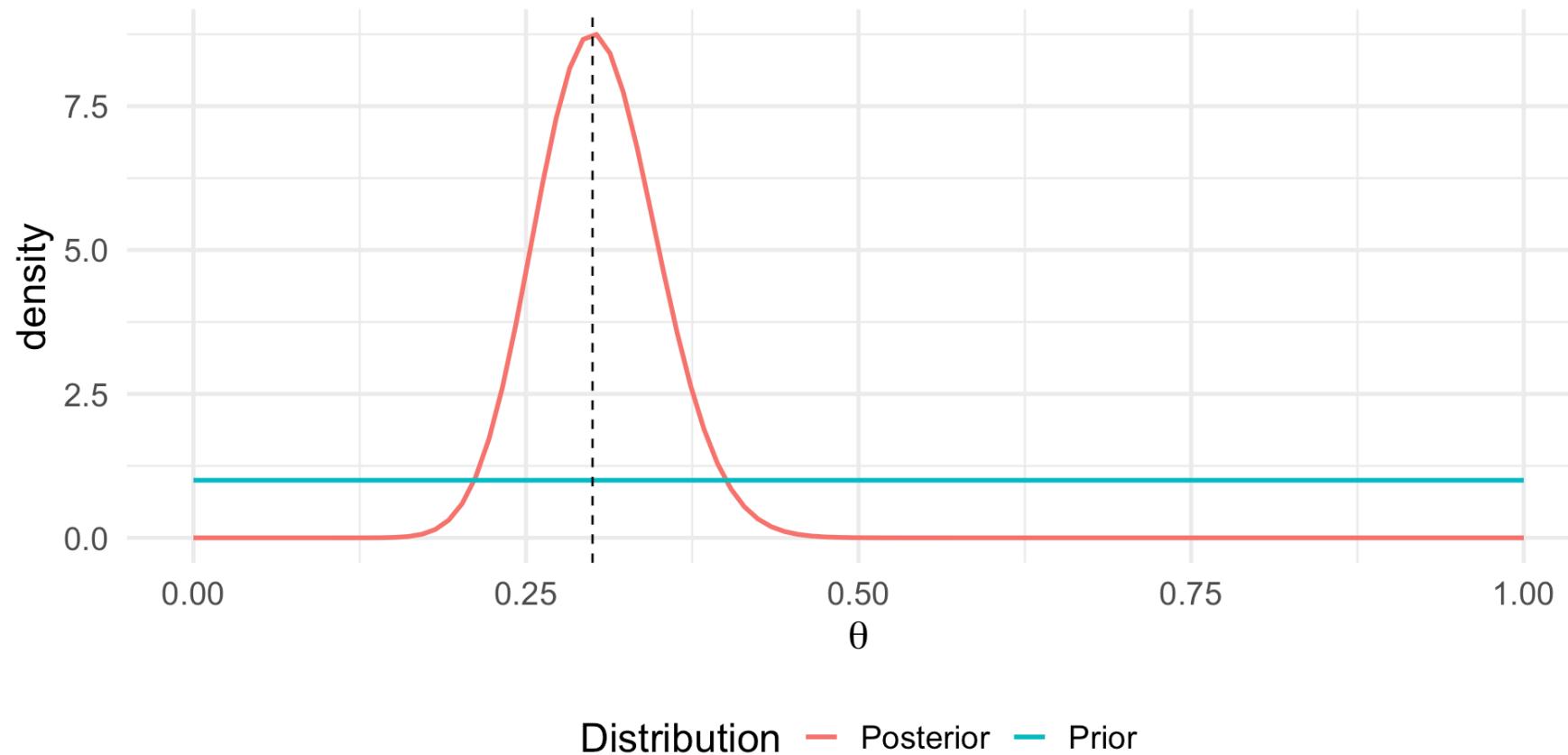
$$\begin{aligned}\pi(\theta|d) &\propto \pi(\theta) f(d|\theta) \\ &\propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^d (1-\theta)^{n-d} \\ &\propto \theta^{\alpha+d-1} (1-\theta)^{\beta+n-d-1}\end{aligned}$$

Numerical example:  $\alpha = \beta = 1, n = 10, d = 3$



$$\mathbb{E}(\Theta|D=d) = \frac{4}{12} = \frac{1}{3}, \quad \text{SD}(\Theta|D=d) = \sqrt{\frac{4 \times 8}{12^2 \times 13}} \approx 0.1307$$

2nd numerical example:  $\alpha = \beta = 1, n = 100, d = 30$



$$\mathbb{E}(\Theta|D = d) = \frac{31}{102} \approx 0.3039, \quad \text{SD}(\Theta|D = d) = \sqrt{\frac{31 \times 71}{102^2 \times 103}} \approx 0.0453$$

Consider now that we want to estimate  $\psi = g(\theta)$ , a coordinate of  $\theta$  or a determinist function of  $\theta$ :

## Credible intervals

- A credible interval of posterior probability  $(1 - \alpha)$  is an interval  $[a, b]$  such that

$$\mathbb{P}(a \leq g(\Theta) \leq b | D = d) = 1 - \alpha.$$

- Quantile credible intervals are computed as:
  - $a$  is the  $\alpha/2$  quantile of the posterior distribution  $[g(\Theta) | D = d]$
  - $b$  is the  $1 - \alpha/2$  quantile of the posterior distribution  $[g(\Theta) | D = d]$

- We have a strong belief that  $\psi = g(\theta)$  lies in this interval.
- There is a real difference with the definition of a confidence interval in the frequentist paradigm

# 1.4 The problems

Given a statistical model  $D \sim P_\theta$ , we want

1. In the frequentist paradigm: to approximate the distribution of  $s(D, \lambda)$  when  $D \sim P_\theta$  for all  $\theta$ 
  - asymptotic approximations for the MLE
  - the bootstrap
2. In the Bayesian paradigm: to compute the posterior distribution of  $[\Theta | D = d]$ 
  - MCMC methods
  - variational inference (not for this course)

# 2. The frequentist paradigm

In this part, we will assume that  $\theta \in \mathbb{R}^d$

## 2.1 Fisher and the maximum likelihood estimator (MLE)

- The **likelihood** of  $\theta$  is the density of  $[D|\Theta = \theta]$  wrt. a dominating measure:  $f(d|\theta)$
- The **log-likelihood** is the logarithm of the likelihood:  $\ell(d|\theta) = \log f(d|\theta)$
- The **maximum likelihood estimator** (MLE) is

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \mathbb{R}^d} \ell(d|\theta)$$

- The **score function** is the gradient of the log-likelihood:

$$\nabla_{\theta} \ell(d|\theta) = \frac{1}{f(d|\theta)} \nabla_{\theta} f(d|\theta)$$

Now, let us consider the random dataset  $D \sim P_\theta$  and the rv  $\widehat{\Theta}$  that models the MLE

- The MLE cancels the score function:  $\nabla_\theta \ell(D|\widehat{\Theta}) = 0$
- Whatever the true value  $\theta$ , we have

$$\begin{aligned}\mathbb{E}_\theta \left[ \nabla_\theta \ell(D|\widehat{\Theta}) \right] &= \int \nabla_\theta \ell(d|\theta) f(d|\theta) dd = \int \frac{1}{f(d|\theta)} \nabla_\theta f(d|\theta) f(d|\theta) dd \\ &= \int \nabla_\theta f(d|\theta) dd = \nabla_\theta \left\{ \int f(d|\theta) dd \right\} = \nabla_\theta 1 = 0\end{aligned}$$

- I.e., at the true value  $\theta$ , the score function has expectation zero. How far is the score function from zero?
- **Fisher information:**  $\mathcal{I}_\theta$  is the variance matrix:

$$\mathcal{I}_\theta = \mathbb{V}_\theta \left[ \nabla_\theta \ell(D|\Theta) \right] = -\mathbb{E}_\theta \left[ \nabla_{\theta\theta}^2 \ell(D|\Theta) \right]$$

- Hence, if we draw a random  $D \sim P_\theta$ , the random score function  $\nabla_\theta \ell(D|\theta)$  at the true  $\theta$  has
  - expectation zero and
  - variance  $\mathcal{I}_\theta$
- Note that  $\mathcal{I}_\theta$  does not depend on the data at hand!
- $\mathcal{I}_\theta$  measures
  - the **local curvature** of the log-likelihood at  $\theta$
  - the amount of information from any random dataset  $D \sim P_\theta$  about  $\theta$

- With a Taylor approx fo the score function at  $\theta$ , we have

$$0 = \nabla_{\theta} \ell(d|\hat{\theta}) \approx \nabla_{\theta} \ell(d|\theta) + \mathcal{J}(d|\theta)(\hat{\theta} - \theta), \quad \text{where } \mathcal{J}(d|\theta) = -\nabla_{\theta\theta}^2 \ell(d|\theta)$$

- Hence,  $\hat{\theta} \approx \theta + \mathcal{J}^{-1}(d|\theta) \nabla_{\theta} \ell(d|\theta) \approx \theta + \mathcal{I}_{\theta}^{-1} \nabla_{\theta} \ell(d|\theta)$
- The approx. error term  $\mathcal{I}_{\theta}^{-1} \nabla_{\theta} \ell(d|\theta)$ 
  - has expectation zero and
  - variance  $\mathcal{I}_{\theta}^{-1} \mathcal{I}_{\theta} \mathcal{I}_{\theta}^{-1} = \mathcal{I}_{\theta}^{-1}$
- Hence Fisher Gaussian approximation:  $\widehat{\Theta} \sim \mathcal{N}(\theta, \mathcal{I}_{\theta}^{-1})$
- Since  $\mathcal{I}_{\theta} = -\mathbb{E}_{\theta} \left[ \nabla_{\theta\theta}^2 \ell(D|\theta) \right]$ , we can approximate  $\mathcal{I}_{\theta}$  by the **observed information matrix**  $\mathcal{J}(d|\theta) = -\nabla_{\theta\theta}^2 \ell(d|\theta)$

# Numerically

- The score function  $\nabla_{\theta}\ell(d|\theta)$  can be computed by **automatic differentiation**, as well as the observed information matrix

$$\mathcal{J}(d|\theta) = -\nabla_{\theta\theta}^2\ell(d|\theta)$$



- When the data  $D = (D_1, \dots, D_n)$  form an iid sample of size  $n$ ,  $P_{\theta} = \bigotimes_{i=1}^n p_{\theta}$  and

$$\mathcal{I}_{\theta} = ni_{\theta}, \quad \text{where } i_{\theta} = -\mathbb{E}_{\theta} [\nabla_{\theta\theta}^2\ell(D_1|\theta)]$$

- The last expected value can be approximated by **Monte Carlo methods...**

- In this case, a specific MC approximation is  $\frac{1}{n} \mathcal{J}(d|\theta) = -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta\theta}^2\ell(d_i|\theta)$ .

Thus  $\mathcal{J}(d|\theta)$  is a **Monte Carlo** approximation of  $\mathcal{I}_{\theta}$

## The binomial case : $D \sim \text{Binomial}(n, \theta)$

```

1 import jax.numpy as jnp
2 from jax import grad, hessian
3
4 # Define the log-likelihood of the binomial distribution
5 def log_likelihood(theta, data, n):
6     return jnp.sum(data*jnp.log(theta) + (n-data)*jnp.log(1-theta))
7
8 # Define the score function and the Hessian
9 score = grad(log_likelihood, argnums=0)
10 hess_loglik = hessian(log_likelihood, argnums=0)
11
12 # The observed data: 30 successes out of 100 trials
13 data, n = 30, 100
14 mle = data/n
15 mle_score = score(mle, data, n)
16 observed_info = - hess_loglik(mle, data, n)
17 SErr = jnp.sqrt(1/observed_info)
18 print(f"The MLE is {mle:.4f}, its score is {mle_score:.4f} and the standard error is {SErr:.4f}")

```

The MLE is 0.3000, its score is -0.0000 and the standard error is 0.0458.

The Gaussian case:  $D \sim \mathcal{N}^{\otimes 100}(\mu, \sigma^2)$ ,  $\theta = (\mu, \sigma^2)$

```

1 # Define the log-likelihood of the Gaussian distribution
2 def log_likelihood_gaussian(theta, data):
3     mu, sigma2 = theta[0], theta[1]
4     return -0.5*jnp.sum(jnp.log(2*jnp.pi*sigma2) + (data-mu)**2/sigma2)
5 # Define the score function and the Hessian
6 score_gaussian = grad(log_likelihood_gaussian, argnums=0)
7 hess_loglik_gaussian = hessian(log_likelihood_gaussian, argnums=0)
8 # The observed data: sample of size 100 from a N(0,1)
9 import numpy as np
10 np.random.seed(1234)
11 data = jnp.array(np.random.normal(0, 1, 100))
12 # MLE
13 mle = jnp.array([jnp.mean(data), jnp.var(data)])
14 mle_score = score_gaussian(mle, data)
15 observed_info = - hess_loglik_gaussian(mle, data)
16 SErr_t, SErr_s = jnp.sqrt(1/observed_info[0, 0]), jnp.sqrt(1/observed_info[1, 1])
17 print(f"The MLE is ({mle[0]:.4f}, {mle[1]:.4f}), its score is ({mle_score[0]:.4f}, {mle_score[1]:.4f})")

```

The MLE is (0.0351, 0.9914), its score is (0.0000, 0.0000) and the standard errors are 0.0996, 0.1402.

```
1 print(f"The observed information matrix is\n{jnp.round(observed_info, decimals = 4)}")
```

The observed information matrix is

```
[[100.8686  0.      ]
 [ 0.       50.8724]].
```

## 2.2 The bootstrap

- The **bootstrap** is a **non-parametric** method to approximate the distribution of a statistic  $s(D, \lambda)$ 
  - where  $D \sim P_\theta$  is the random dataset  $\sim P_\theta$
  - and  $\lambda$  is the tuning parameter
- It can be applied to any statistic  $s(D, \lambda)$ , not only the MLE
- The dataset  $D$  should be composed of  $(D_1, \dots, D_n)$  iid rv's, with  $D_i \sim p_\theta$ , i.e.

$$P_\theta = \bigotimes_{i=1}^n p_\theta$$

- In the above notations:
  - $p_\theta$  is the marginal **distribution of one individual**  $D_i$
  - $P_\theta$  is the **joint distribution of the dataset**  $D$ , which is a product (because of independence)

- The idea of bootstrap is to approximate the marginal  $p_\theta$  by the empirical distribution of the data, as in Monte Carlo methods:

$$p^* = \frac{1}{n} \sum_{i=1}^n \delta_{d_i}$$

- Thus,  $p_\theta \approx p^*$  and  $P_\theta \approx P^* = \bigotimes_{j=1}^n p^*$  (the population is replaced by the sample)

## The bootstrap principle

- The distribution of  $s(D, \lambda)$  when  $D \sim P_\theta$  can be approximated by the distribution of  $s(D^*, \lambda)$  when  $D^* \sim P^*$
- It is easy to draw  $D^* \sim P^*$  by **drawing  $n$  times with replacement** from the data  $D$
- A dataset  $d^*$  drawn from  $P^*$  is called a **bootstrap sample**
- It has the same size  $n$  as the original dataset  $d$

## Example with $\mathcal{N}^{\otimes 100}(\mu, \sigma^2)$ to compute the bootstrap standard errors

```

1 def estimate(data, lmbda=0):
2     return np.array([np.mean(data), np.var(data, ddof=lmbda)])
3 data2 = np.array(data)
4 # Bootstrap
5 N = 10000
6 bootstrap_samples = \
7     np.array([estimate(np.random.choice(data2, size=100, replace=True), 0) for _ in ran
8 bootstrap_var = np.var(bootstrap_samples, axis=0)
9 print(f"The bootstrap standard errors are {np.sqrt(bootstrap_var[0]):.4f} and {np.sqrt(

```

The bootstrap standard errors are 0.0994 and 0.1729.

```

1 print(f"The standard errors obtained with the normal approximation were {SErr_t:.4f} an

```

The standard errors obtained with the normal approximation were 0.0996 and 0.1402.

We can also get Standard errors with  $\lambda = 1$  (unbiased estimator of the variance)

```
1 # Bootstrap
2 bootstrap_samples = \
3     np.array([estimate(np.random.choice(data2, size=100, replace=True), 1) for _ in ran
4 bootstrap_var = np.var(bootstrap_samples, axis=0)
5 print(f"The bootstrap standard errors are now {np.sqrt(bootstrap_var[0]):.4f} and {np.s
```

The bootstrap standard errors are now 0.0994 and 0.1762.

- It can be applied to any statistic, not only the MLE
- The **bootstrap** is a **non-parametric** method, i.e., no assumption on the distribution of the data
- The **bootstrap** is a **resampling** method, i.e., it draws  $n$  times with replacement from the data
- The **bootstrap** is particularly useful to approximate the **standard errors** of a statistic  $s(D, \lambda)$ :

$$\widehat{\mathbb{V}_\theta(s(D, \lambda))} = \mathbb{V}^*(s(D^*, \lambda)), \quad \text{where } D^* \sim P^*$$

The above variance being itself approximated with a basic Monte Carlo method

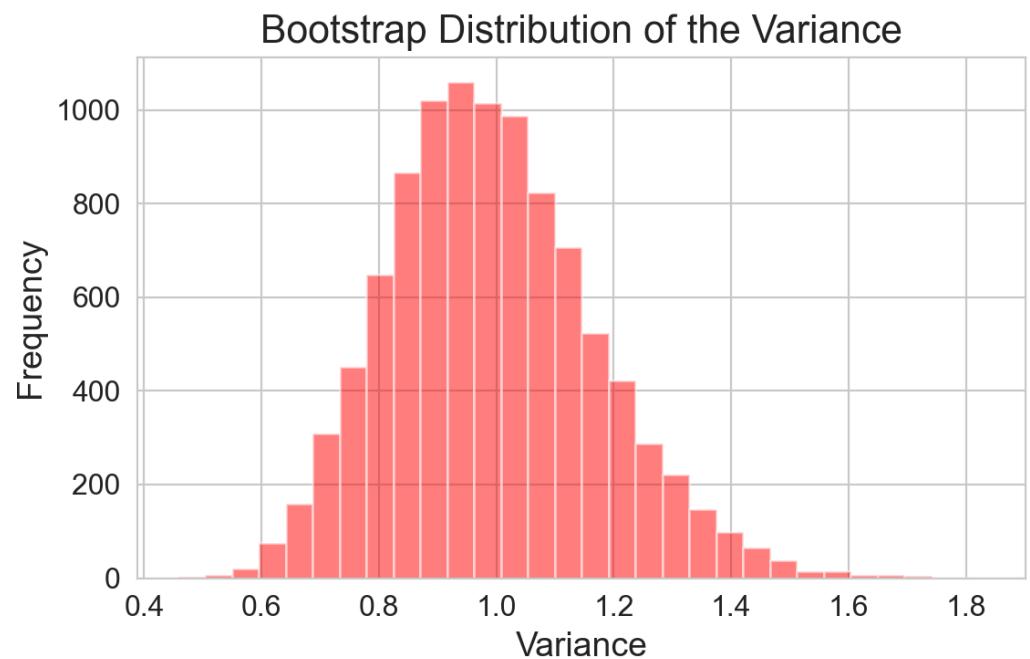
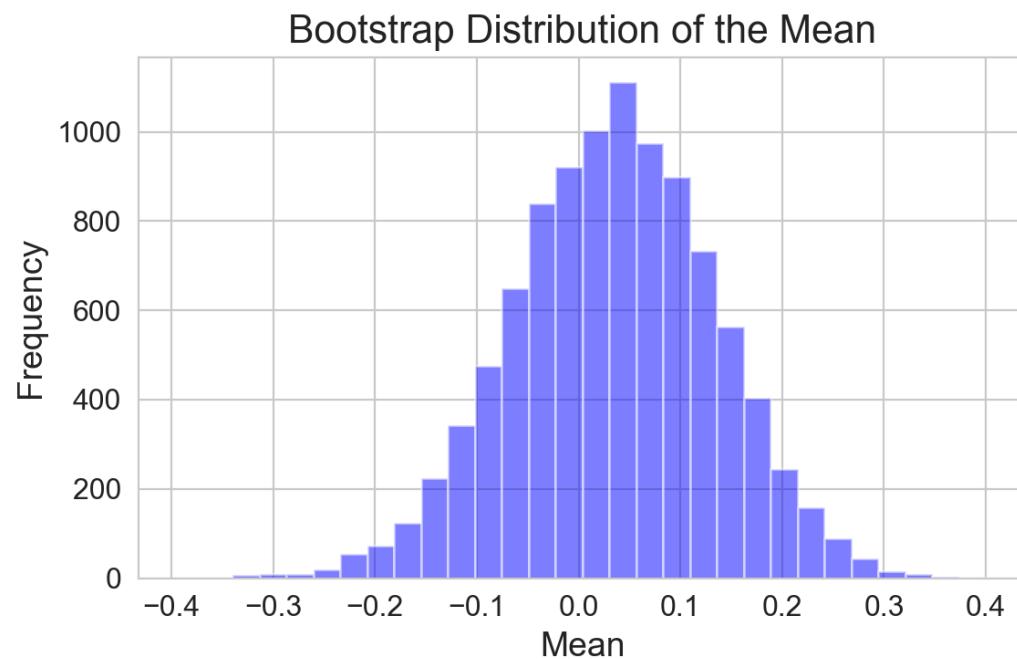
# Confidence intervals with the bootstrap

Consider now that we want to estimate  $\psi$ , and a statistic  $s(D, \lambda)$  that estimates  $\psi$ :

- $\psi$  is a given coordinate of  $\theta$
- or  $\psi = \beta_0 + \sum_j \beta_j x_{n+1,j}$  in (generalized) linear models
- The **percentile method** to get a  $1 - \alpha$  confidence interval for  $\psi$  is obtained as follows:
  - draw  $B$  bootstrap samples  $d_1^*, \dots, d_B^*$  from  $P^*$
  - compute the  $B$  bootstrap estimates  $\hat{\psi}_1^*, \dots, \hat{\psi}_B^*$  of  $\psi$
  - compute the quantiles of order  $\alpha/2$  and  $1 - \alpha/2$  of the bootstrap estimates
  - the percentile confidence interval of level  $(1 - \alpha)$  is  $[\hat{\psi}_{(\alpha/2)}^*, \hat{\psi}_{(1-\alpha/2)}^*]$

- To understand the percentile method, we can
  - consider the bootstrap estimates as a sample drawn a Bayesian posterior distribution
  - approximate a credible interval of posterior probability  $(1 - \alpha)$  with the quantiles of order  $\alpha/2$  and  $1 - \alpha/2$
  - consider that the credible interval is an approximation of a confidence interval
- The percentile method is a **non-parametric** method, i.e., no assumption on the distribution of the data
- The percentile method is a **first** approximation of the confidence interval, but it is easy to implement
- **Better approximations** exist, e.g., the **BCa method** (bias-corrected and accelerated), but will not be studied in this course

## Example with the Gaussian case (same data as before)



```

1 # Quantiles
2 alpha = 0.05
3 quantiles = np.quantile(bootstrap_samples, [alpha/2, 1-alpha/2], axis=0)
4 print(f"The percentile confidence interval is [{quantiles[0, 0]:.4f}, {quantiles[1, 0]:.4f}]")

```

The percentile confidence interval is [-0.1624, 0.2253] for the mean.

```

1 print(f"The percentile confidence interval is [{quantiles[0, 1]:.4f}, {quantiles[1, 1]:.4f}]")

```

The percentile confidence interval is [0.6845, 1.3744] for the variance.

# Is the bootstrap always working?

- The answer is definitely **no!**
- E.g., if the estimator is constant,  $s(D, \lambda) = \check{\theta}$  for a given  $\check{\theta}$  whatever the data  $D$ , the bootstrap will not work
  - the bootstrap samples will be all equal to  $\check{\theta}$
  - the bootstrap standard errors will be zero
- Remember that, for an estimator  $\widehat{\Theta} = s(D, \lambda)$ , the **mean square error** is

$$\text{MSE}(\widehat{\Theta}) = \mathbb{E}_\theta \left( (\widehat{\Theta} - \theta)^2 \right) = \mathbb{V}_\theta(\widehat{\Theta}) + \text{Bias}_\theta^2(\widehat{\Theta})$$

- Bootstrap methods cannot estimate the bias
- Thus, the bootstrap requires that

$$\text{Bias}_{\theta}^2(\widehat{\Theta}) \ll \mathbb{V}_{\theta}(\widehat{\Theta})$$

- The MLE of the parameter of an iid model  $P_{\theta} = \otimes_{i=1}^n p_{\theta}$  satisfies this condition if the size  $n$  is large:
  - the bias is of order  $1/n$  (if not zero)
  - the variance is of order  $1/n$
- Any unbiased estimator satisfies this condition...

# 3. The Bayesian paradigm

Now,  $\Theta$  is a rv, with marginal probability density function  $\pi(\theta)$ , which is the **prior density**

# 3.1 The posterior distribution

- In Bayesian statistics, the **likelihood**  $f(d|\theta)$  is used only as a **weight** to update the prior  $\pi(\theta)$  regarding its fit to the data  $d$
- The **posterior density** of  $\Theta$  given the data  $d$  has density

$$\pi(\theta|d) = \frac{f(d|\theta)\pi(\theta)}{f_\pi(d)} \propto f(d|\theta)\pi(\theta),$$

where  $f_\pi(d)$  is the normalizing constant.

- Apart from very specific cases, we cannot recognize a well known distribution for the posterior law
- Instead, we used **Monte Carlo methods to get a sample**  $\theta_1, \dots, \theta_N$  from the posterior

## 3.2 Metropolis-Hastings algorithm

- We have already met a distribution known up to a multiplicative constant: the Gibbs measures in simulated annealing
- The MCMC method that was introduced then was a **Metropolis-Hastings algorithm**
- The Metropolis-Hastings algorithm is a **Markov chain Monte Carlo (MCMC) method** that generates a sample from a distribution known up to a multiplicative constant
- Let us denote  $\pi_u(\theta|d)$  an **unnormalized version** of the posterior density. E.g.  
$$\pi_u(\theta|d) = f(d|\theta)\pi(\theta)$$
- (We can get rid of any multiplicative term that does not depend on  $\theta$  in  $\pi_u(\theta|d)$ )

## Metropolis-Hastings Algorithm

1. Initialize  $\theta_0$
2. For  $t = 1, 2, \dots, T$ :
  - Generate  $\zeta \sim g(\zeta|\theta_t)$ , a proposal kernel given  $\theta_t$
  - Compute the acceptance ratio  $\rho = \frac{\pi_u(\zeta|d)g(\theta_t|\zeta)}{\pi_u(\theta_t|d)g(\zeta|\theta_t)}$
  - Draw  $U \sim \text{Unif}(0, 1)$  and set  $\theta_{t+1} = \zeta \mathbf{1}\{U < \rho\} + \theta_t \mathbf{1}\{U \geq \rho\}$

- The **proposal kernel**  $g(\zeta|\theta)$  is a density that generates a candidate  $\zeta$  given the current value  $\theta$

The same remarks as for the simulated annealing algorithm apply here:

- if  $g$  is symmetric, the acceptance ratio simplifies
- **choose the proposal kernel  $g(\zeta|\theta)$  wisely**
- the average proportion of accepted moves along time should be around 0.234
- the **burn-in period** is the time needed to reach the stationary distribution
  - assume we start from a  $\theta_0$  from the prior
  - we can get rid of the first  $B$  iterations until we reach the part of the sample where the posterior distribution is not near 0.

- We can also start from the  $\operatorname{argmax}_{\theta} \pi_u(\theta|d)$ 
  - this is the **maximum a posteriori (MAP)** estimator
  - it can often be computed numerically with an optimization algorithm...
- During the burn-in period, we can do anything we want, in particular change the proposal kernel
- At the end, the posterior distribution is approximated by

$$\pi(\theta|d)d\theta \approx \frac{1}{T-B} \sum_{t=B+1}^T \delta_{\theta_t}$$

# Numerical example

- Statistical model:  $D = (D_1, \dots, D_n) \sim \mathcal{N}^{\otimes n}(\mu, \sigma^2)$
- Prior:  $\pi(\mu, \sigma^2) = \pi(\mu|\sigma^2)\pi(\sigma^2)$ , where  $[\mu|\sigma^2] \sim \mathcal{N}(0, \sigma^2)$  and  $1/\sigma^2 \sim \text{Exp}(1)$
- Thus,  $\pi(\sigma^2)d\sigma^2 = \exp(-\sigma^{-2})\sigma^{-4}d\sigma^2$
- The likelihood is

$$\begin{aligned}
 f(d|\theta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(d_i - \mu)^2}{2\sigma^2}\right) \propto \frac{1}{\sigma^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (d_i - \mu)^2\right) \\
 &\propto \frac{1}{\sigma^n} \exp\left(-\frac{n\mu^2}{2\sigma^2} - \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^n d_i^2}_{B} + \underbrace{\frac{\mu}{\sigma^2} \sum_{i=1}^n d_i}_{A}\right)
 \end{aligned}$$

depends on the data only through  $A$  and  $B$

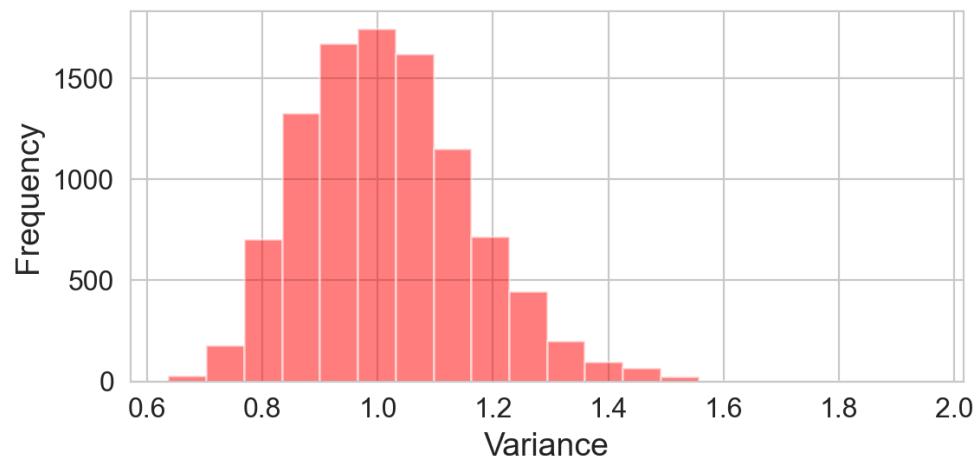
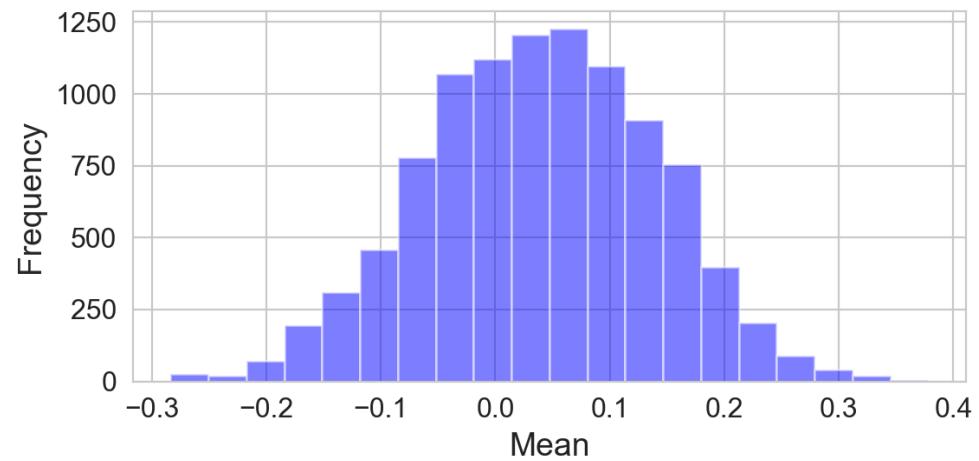
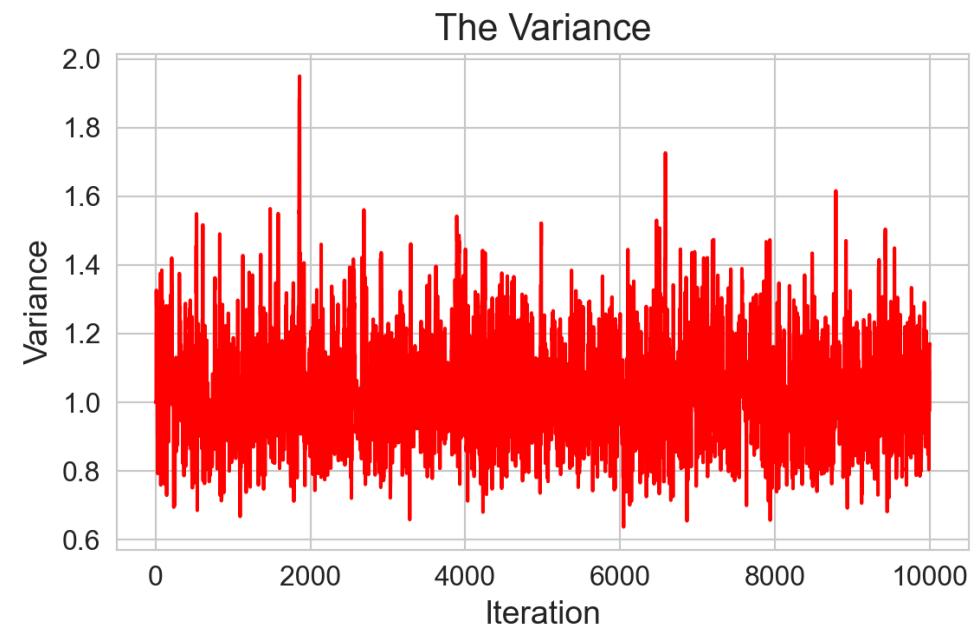
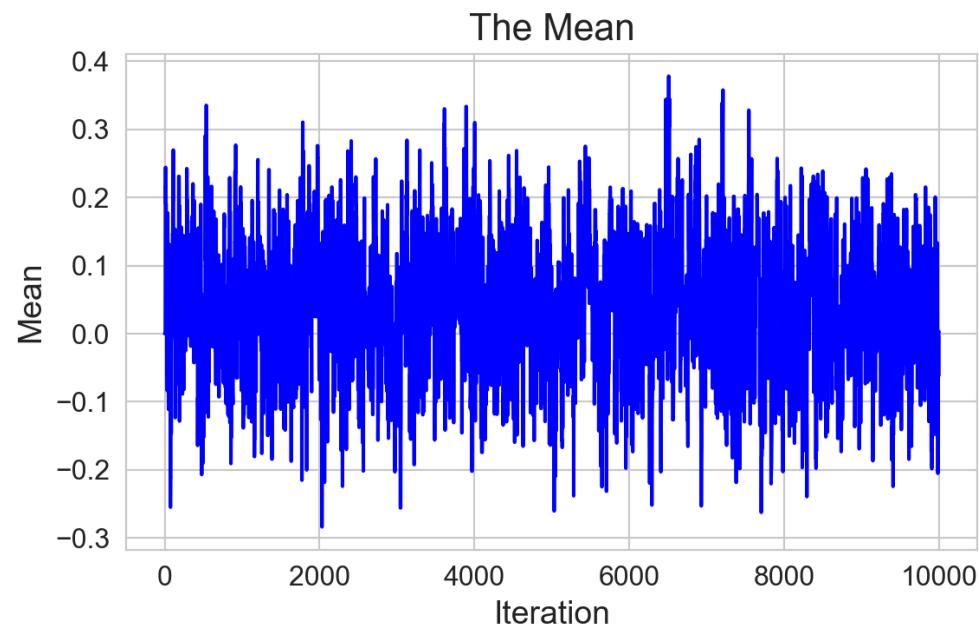
Always the same data of size 100, drawn from a  $\mathcal{N}(0, 1)$

```

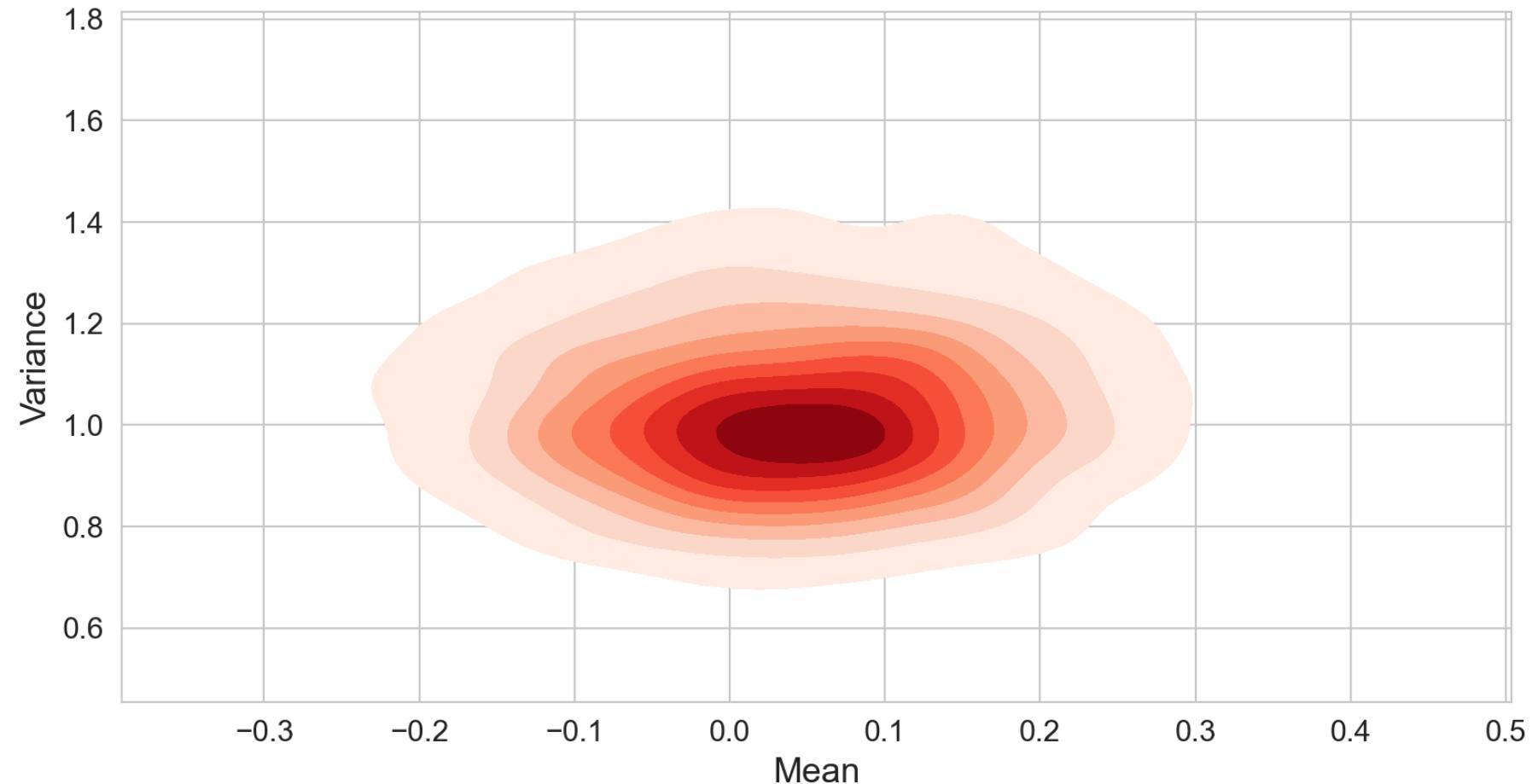
1 import numpy as np
2 from scipy.stats import norm, expon, multivariate_normal
3 A, B = np.sum(data2), np.sum(data2**2)
4 def log_posterior(theta, n=100):
5     if theta[1] <= 0:
6         return -np.inf
7     else:
8         mu, sigma2 = theta[0], theta[1]
9         return -1./sigma2 - 2.*np.log(sigma2) - 0.5*n*np.log(sigma2) - n*mu**2/sigma2/2 +
10
11 # Metropolis-Hastings
12 def metropolis_hastings(log_posterior, T, theta0, g):
13     theta = theta0
14     samples = [theta]
15     for t in range(T):
16         zeta = theta + g.rvs(size = 1)
17         rho = np.exp(log_posterior(zeta) - log_posterior(theta))
18         U = np.random.uniform()
19         theta = zeta if U < rho else theta
20         samples.append(theta)
21     return np.array(samples)

```

## Paths of the algo. and posterior distributions



## Bivariate heatmap of the posterior samples



## Summaries of the posterior distribution

```

1 # Summaries of the posterior distribution
2 mean_posterior = np.mean(posterior_sample, axis=0)
3 median_posterior = np.median(posterior_sample, axis=0)
4 cov_posterior = np.round(np.cov(posterior_sample, rowvar=False), decimals=4)
5 print(f"The posterior mean is {mean_posterior[0]:.4f} for the mean and {mean_posterior[

```

The posterior mean is 0.0399 for the mean and 1.0215 for the variance.

```

1 print(f"The posterior covariance matrix is\n{cov_posterior}.")

```

The posterior covariance matrix is  

$$\begin{bmatrix} 0.0102 & 0.0002 \\ 0.0002 & 0.0203 \end{bmatrix}.$$

```

1 print(f"The posterior median is {median_posterior[0]:.4f} for the mean and {median_post

```

The posterior median is 0.0417 for the mean and 1.0046 for the variance.

## 3.3 The Gibbs algorithm

- The Gibbs algorithm relies on a **decomposition of the parameter**  $\theta$  into  $K$  blocks

$$\theta = (\theta^{(1)}, \dots, \theta^{(K)}).$$

- We assume that we are able to compute the **full conditional distributions** of each block  $\theta^{(k)}$  given the other blocks  $\theta^{(-k)}$  and the data  $d$

$$\pi(\theta^{(k)}|d, \theta^{(-k)}) = \frac{\pi(\theta^{(k)}, \theta^{(-k)}|d)}{\pi(\theta^{(-k)}|d)} \propto \pi(\theta|d)$$

- The Gibbs algorithm updates one block given the other blocks and the data to obtain the new value  $\theta_{t+1}$ .

## Gibbs Algorithm with random block selection

1. Initialize  $\theta_0$
2. For  $t = 1, 2, \dots, T$ :
  - Draw  $k \sim \text{Unif}(1, K)$
  - Draw  $\theta_{t+1}^{(k)} \sim \pi(\theta^{(k)} | d, \theta_t^{(-k)})$  and set  $\theta_{t+1}^{(-k)} = \theta_t^{(-k)}$

## Gibbs Algorithm with systematic block selection

1. Initialize  $\theta_0$
2. For  $t = 1, 2, \dots, T$ :
  - Set  $k = (t \mod K) + 1$
  - Draw  $\theta_{t+1}^{(k)} \sim \pi(\theta^{(k)} | d, \theta_t^{(-k)})$  and set  $\theta_{t+1}^{(-k)} = \theta_t^{(-k)}$

- Like the Metropolis-Hastings algorithm, the Gibbs algorithm generates a correlated sample from the posterior distribution
- The path  $\theta_0, \dots, \theta_T$  is a Markov chain
- As with Metropolis-Hastings, the Gibbs algorithm requires a **burn-in period**
- At the end, the posterior distribution is approximated by

$$\pi(\theta|d)d\theta \approx \frac{1}{T - B} \sum_{t=B+1}^T \delta_{\theta_t}$$

where  $B$  is the number of burn-in iterations

# Numerical example

- Statistical model:  $D = (D_1, \dots, D_n) \sim \mathcal{N}^{\otimes n}(\mu, \sigma^2)$ , and we set  $\tau = 1/\sigma^2$  and  $\theta = (\mu, \tau)$
- Prior:  $\pi(\mu, \tau) = \pi(\mu|\tau)\pi(\tau)$ , where  $\mu|\tau \sim \mathcal{N}(0, 1/\tau)$  and  $\tau \sim \text{Exp}(1)$
- Likelihood:

$$f(d|\theta) \propto \tau^{n/2} \exp(-n\mu^2\tau/2 - B\tau/2 + \mu\tau A)$$

with  $A = \sum_{i=1}^n d_i$  and  $B = \sum_{i=1}^n d_i^2$

- Posterior:

$$\pi(\mu, \tau|d) \propto \tau^{n/2} \exp(-n\mu^2\tau/2 - B\tau/2 + \mu\tau A) e^{-\tau} \exp(-\tau\mu^2/2)$$

- The full conditional distributions are

- For  $[\mu|\tau, d]$ :

$$\begin{aligned}\pi(\mu|\tau, d) &\propto \exp(-n\mu^2\tau/2 + \mu\tau A - \mu^2/2) \\ &\propto \exp\left\{-\frac{(n+1)\tau}{2}\left(\mu - \frac{A}{n+1}\right)^2\right\}\end{aligned}$$

i.e.,  $[\mu|\tau, d] \sim \mathcal{N}\left(\frac{A}{n+1}, \frac{1}{(n+1)\tau}\right)$

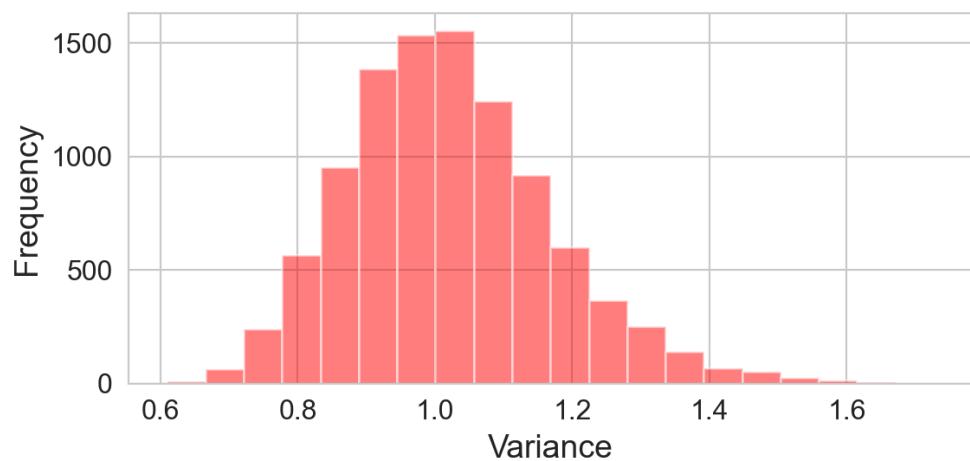
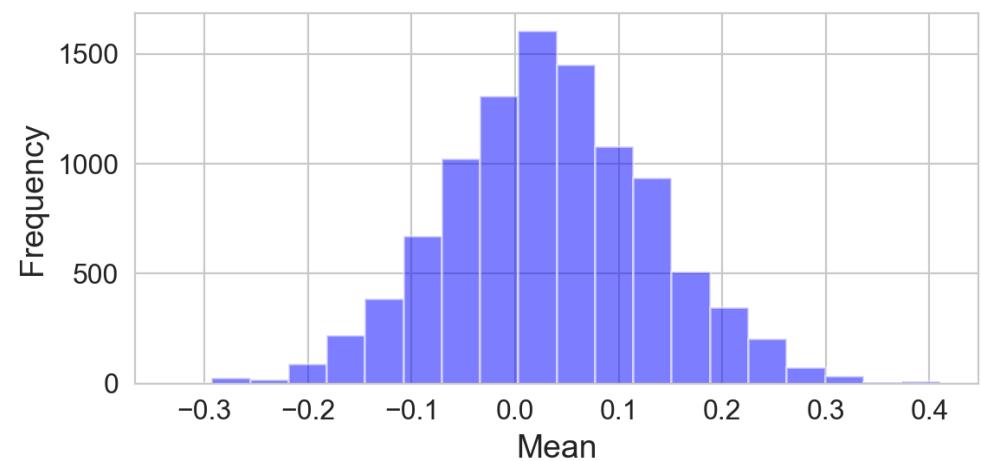
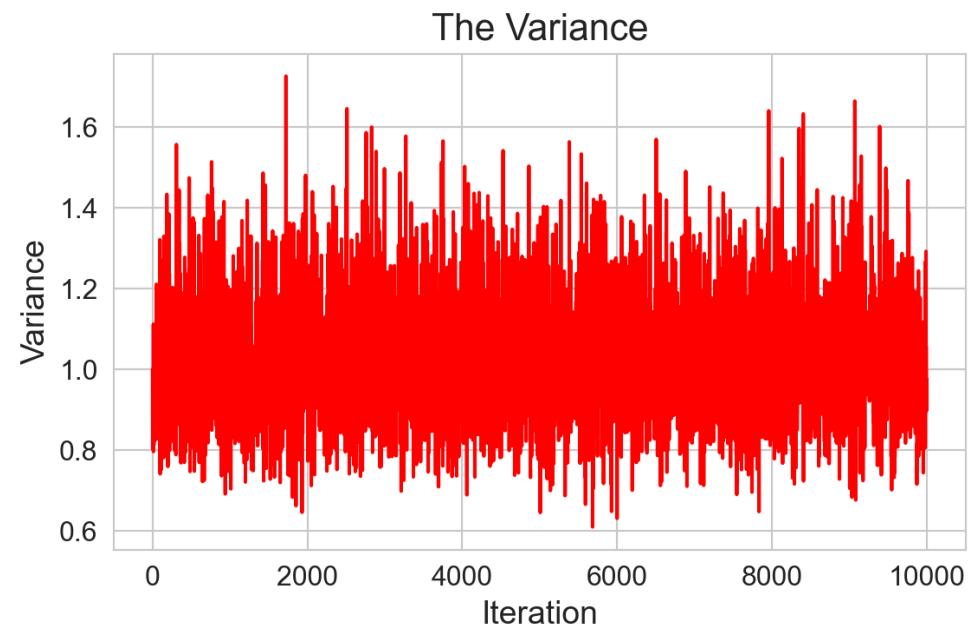
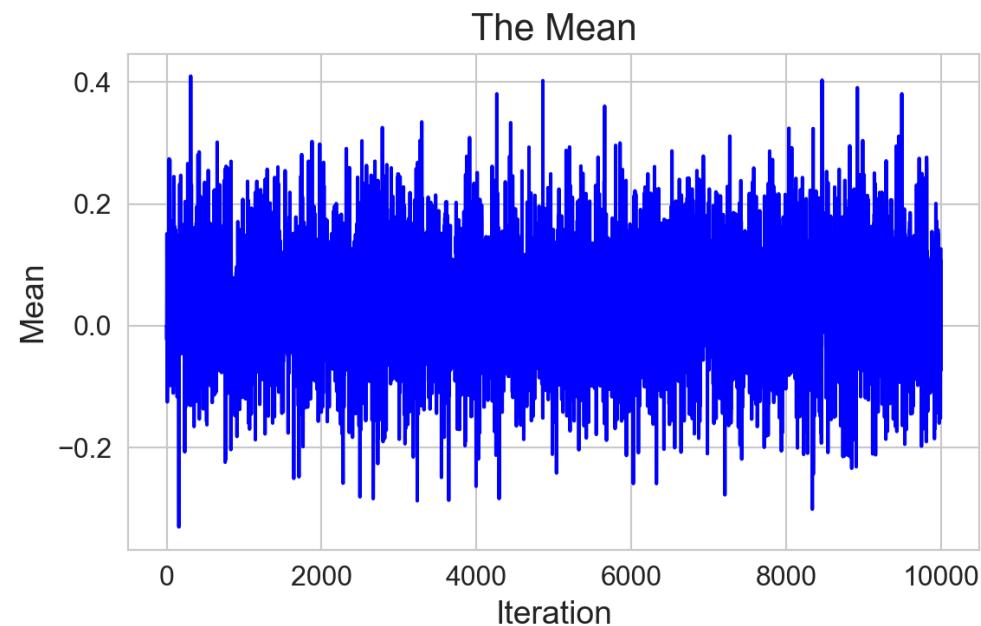
- For  $[\tau|\mu, d]$ :

$$\begin{aligned}\pi(\tau|\mu, d) &\propto \tau^{n/2} \exp(-n\mu^2\tau/2 - B\tau/2 + \mu\tau A) e^{-\tau} \exp(-\tau\mu^2/2) \\ &\propto \tau^{n/2} \exp\left\{-\frac{\tau}{2}((n+1)\mu^2 + B - 2\mu A + 2)\right\}\end{aligned}$$

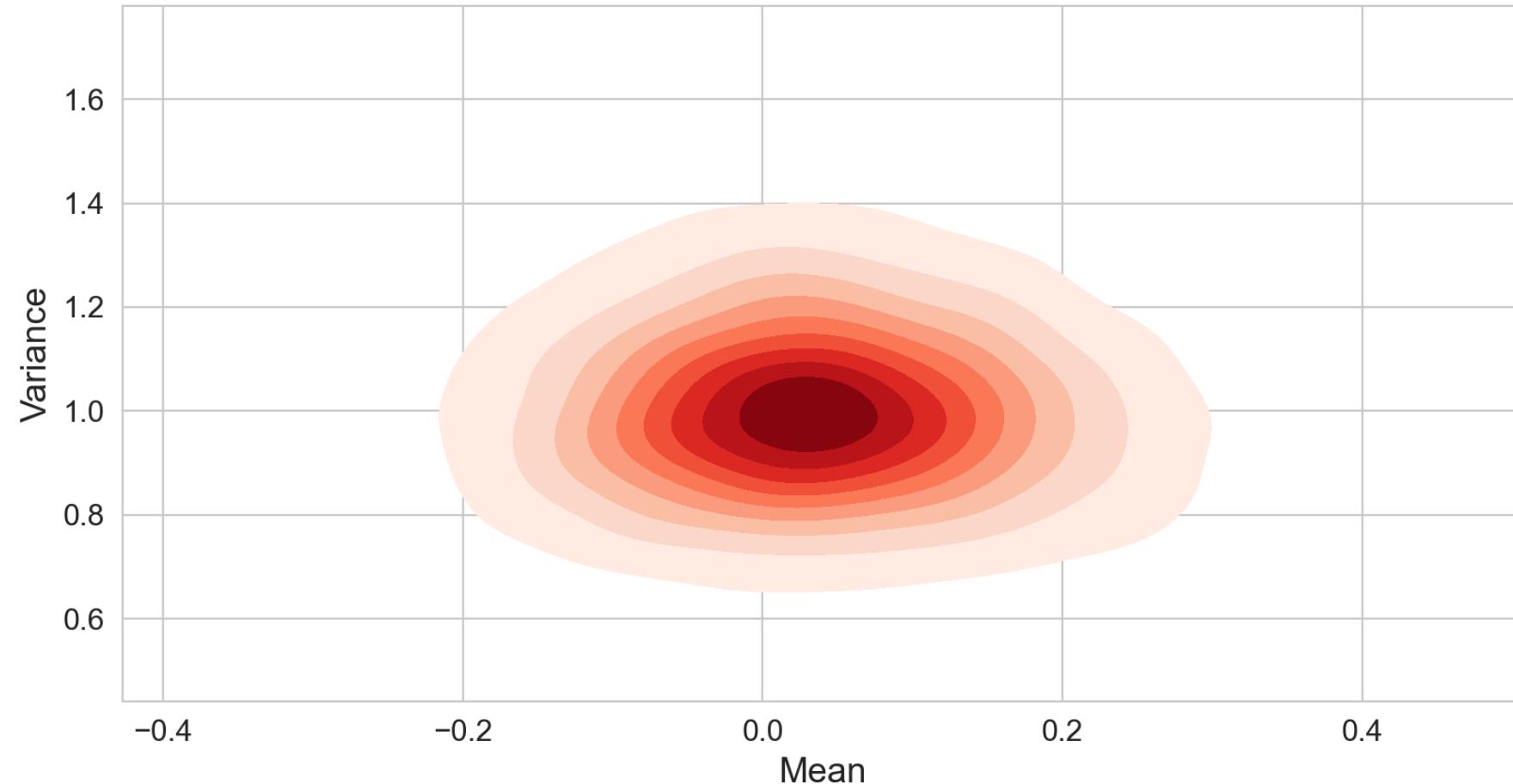
i.e.,  $[\tau|\mu, d] \sim \text{Gamma}\left(\frac{n}{2} + 1, \frac{1}{2}(B - 2\mu A + (n+1)\mu^2) + 1\right)$

```
1 from scipy.stats import norm, gamma
2 def gibbs_sampler(n, data, T, theta0, scan='random'):
3     A, B = np.sum(data), np.sum(data**2)
4     theta = theta0
5     samples = [theta]
6     for t in range(T):
7         if scan == 'random':
8             k = np.random.choice([0, 1])
9         else:
10            k = t % 2
11         if k == 0:
12             mu = np.random.normal(A/(n+1), np.sqrt(1/(n+1)/theta[1]))
13             tau = theta[1]
14         else:
15             mu = theta[0]
16             rate = 0.5*(B - 2*mu*A + (n+1)*mu**2) + 1
17             tau = np.random.gamma(n/2+1, 1/rate)
18             theta = np.array([mu, tau])
19             samples.append(theta)
20     return np.array(samples)
21 # Run
```

## Paths of the algo. and posterior distributions



## Bivariate heatmap of the Gibbs samples



## Summaries of the Gibbs distribution

```

1 # Summaries of the Gibbs distribution
2 mean_gibbs = np.mean(gibbs_sample, axis=0)
3 median_gibbs = np.median(gibbs_sample, axis=0)
4 cov_gibbs = np.round(np.cov(gibbs_sample, rowvar=False), decimals=4)
5 print(f"The Gibbs mean is {mean_gibbs[0]:.4f} for the mean and {mean_gibbs[1]:.4f} for"

```

The Gibbs mean is 0.0343 for the mean and 0.9996 for the variance.

```
1 print(f"The Gibbs covariance matrix is\n{cov_gibbs}.")
```

The Gibbs covariance matrix is  

$$\begin{bmatrix} 0.0099 & -0.0002 \\ -0.0002 & 0.0205 \end{bmatrix}.$$

```
1 print(f"The Gibbs median is {median_gibbs[0]:.4f} for the mean and {median_gibbs[1]:.4f}
```

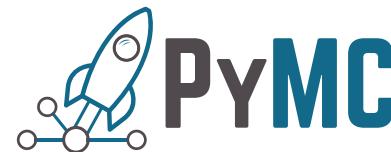
The Gibbs median is 0.0313 for the mean and 0.9917 for the variance.

# Python libraries to perform Bayesian inference

- PyMC: a Python library for probabilistic programming
- Stan: a probabilistic programming language
- Edward: a library for probabilistic modeling, inference, and criticism
- ...

These packages provide high-level languages to define the model and perform the Monte Carlo approximation of the posterior

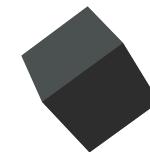
- MCMC: Metropolis-Hastings, Gibbs, Hamiltonian Monte Carlo, No-U-Turn Sampler, ...
- Variational inference (=approximation of the posterior by e.g. a Gaussian, or another family of distribution): Variational Bayes, Black Box Variational Inference, ...



PyMC



Stan



Edward

```
1 import numpy as np  
2 import pymc as pm
```

WARNING (pytensor.tensor.blas): Using NumPy C-API based implementation for BLAS functions.

```
1 import arviz as az  
2 import matplotlib.pyplot as plt  
3  
4 # Modèle bayésien avec pymc  
5 with pm.Model() as model:  
6     tau = pm.Exponential("tau", 1.0) # tau = 1/sigma^2  
7     sigma2 = pm.Deterministic("sigma2", 1 / tau)  
8     mu = pm.Normal("mu", mu=0, sigma=pm.math.sqrt(sigma2))  
9  
10    # Vraisemblance  
11    y_obs = pm.Normal("y_obs", mu=mu, sigma=pm.math.sqrt(sigma2), observed=data2)  
12  
13    # Inférence  
14    trace = pm.sample(2000, cores=1, return_inferencedata=True, progressbar=False)
```

	mean	sd	hdi_3%	hdi_97%	...	mcse_sd	ess_bulk	ess_tail	r_hat
mu	0.036	0.100	-0.148	0.227	...	0.002	4029.0	2925.0	1.0
sigma2	1.010	0.142	0.746	1.266	...	0.002	3755.0	2759.0	1.0

[2 rows x 9 columns]

```
array([<Axes: title={'center': 'mu'}>, <Axes: title={'center': 'sigma2'}>],  
      dtype=object)
```

