

Resampling Techniques, Bootstrap and Blocking

Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no¹

Department of Physics and Center for Computing in Science Education, University
of Oslo, Oslo, Norway¹

March 21, 2025

© 1999-2024, Morten Hjorth-Jensen Email morten.hjorth-jensen@fys.uio.no. Released under CC

Attribution-NonCommercial 4.0 license

Overview of week March 17-21, 2025

Topics

1. Reminder from last week about statistical observables, the central limit theorem and bootstrapping, see notes from last week
2. Resampling Techniques, emphasis on Blocking
3. Discussion of onebody densities (whiteboard notes)
4. Video of lecture TBA

Why resampling methods ?

Statistical analysis

- ▶ Our simulations can be treated as *computer experiments*. This is particularly the case for Monte Carlo methods
- ▶ The results can be analysed with the same statistical tools as we would use analysing experimental data.
- ▶ As in all experiments, we are looking for expectation values and an estimate of how accurate they are, i.e., possible sources for errors.

Statistical analysis

- ▶ As in other experiments, many numerical experiments have two classes of errors:
 1. Statistical errors
 2. Systematical errors
- ▶ Statistical errors can be estimated using standard tools from statistics
- ▶ Systematical errors are method specific and must be treated differently from case to case.

And why do we use such methods?

As you will see below, due to correlations between various measurements, we need to evaluate the so-called covariance in order to establish a proper evaluation of the total variance and the thereby the standard deviation of a given expectation value.

The covariance however, leads to an evaluation of a double sum over the various stochastic variables. This becomes computationally too expensive to evaluate. Methods like the Bootstrap, the Jackknife and/or Blocking allow us to circumvent this problem.

Central limit theorem

Last week we derived the central limit theorem with the following assumptions:

Measurement i

We assumed that each individual measurement x_{ij} is represented by stochastic variables which independent and identically distributed (iid). This defined the sample mean of of experiment i with n samples as

$$\bar{x}_i = \frac{1}{n} \sum_j x_{ij}.$$

and the sample variance

$$\sigma_i^2 = \frac{1}{n} \sum_j (x_{ij} - \bar{x}_i)^2.$$

Further remarks

Note that we use n instead of $n - 1$ in the definition of variance. The sample variance and the sample mean are not necessarily equal to the exact values we would get if we knew the corresponding probability distribution.

Running many measurements

Adding m measurements i

With the assumption that the average measurements i are also defined as iid stochastic variables and have the same probability function p , we defined the total average over m experiments as

$$\bar{X} = \frac{1}{m} \sum_i \bar{x}_i.$$

and the total variance

$$\sigma_m^2 = \frac{1}{m} \sum_i (\bar{x}_i - \bar{X})^2.$$

These are the quantities we used in showing that if the individual mean values are iid stochastic variables, then in the limit $m \rightarrow \infty$, the distribution for \bar{X} is given by a Gaussian distribution with variance σ_m^2 .

Adding more definitions

The total sample variance over the mn measurements is defined as

$$\sigma^2 = \frac{1}{mn} \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{X})^2.$$

We have from the equation for σ_m^2

$$\bar{x}_i - \bar{X} = \frac{1}{n} \sum_{j=1}^n (x_{ij} - \bar{X}),$$

and introducing the centered value $\tilde{x}_{ij} = x_{ij} - \bar{X}$, we can rewrite σ_m^2 as

$$\sigma_m^2 = \frac{1}{m} \sum_i (\bar{x}_i - \bar{X})^2 = \frac{1}{m} \sum_{i=1}^m \left[\frac{1}{n} \sum_{j=1}^n \tilde{x}_{ij} \right]^2.$$

Further rewriting

We can rewrite the latter in terms of a sum over diagonal elements only and another sum which contains the non-diagonal elements

$$\begin{aligned}\sigma_m^2 &= \frac{1}{m} \sum_{i=1}^m \left[\frac{i}{n} \sum_{j=1}^n \tilde{x}_{ij} \right]^2 \\ &= \frac{1}{mn^2} \sum_{i=1}^m \sum_{j=1}^n \tilde{x}_{ij}^2 + \frac{2}{mn^2} \sum_{i=1}^m \sum_{j < k}^n \tilde{x}_{ij} \tilde{x}_{ik}.\end{aligned}$$

The first term on the last rhs is nothing but the total sample variance σ^2 divided by m . The second term represents the covariance.

The covariance term

Using the definition of the total sample variance we have

$$\sigma_m^2 = \frac{\sigma^2}{m} + \frac{2}{mn^2} \sum_{i=1}^m \sum_{j < k}^n \tilde{x}_{ij} \tilde{x}_{ik}.$$

The first term is what we have used till now in order to estimate the standard deviation. However, the second term which gives us a measure of the correlations between different stochastic events, can result in contributions which give rise to a larger standard deviation and variance σ_m^2 . Note also the evaluation of the second term leads to a double sum over all events. If we run a VMC calculation with say 10^9 Monte carlo samples, the latter term would lead to 10^{18} function evaluations. We don't want to, by obvious reasons, to venture into that many evaluations.

Note also that if our stochastic events are iid then the covariance terms is zero.

Rewriting the covariance term

We introduce now a variable $d = |j - k|$ and rewrite

$$\frac{2}{mn^2} \sum_{i=1}^m \sum_{j < k}^n \tilde{x}_{ij} \tilde{x}_{ik},$$

in terms of a function

$$f_d = \frac{2}{mn} \sum_{i=1}^m \sum_{k=1}^{n-d} \tilde{x}_{ik} \tilde{x}_{i(k+d)}.$$

We note that for $d = 0$ we have

$$f_0 = \frac{2}{mn} \sum_{i=1}^m \sum_{k=1}^n \tilde{x}_{ik} \tilde{x}_{i(k)} = \sigma^2!$$

Introducing the correlation function

We introduce then a correlation function $\kappa_d = f_d/\sigma^2$. Note that $\kappa_0 = 1$. We rewrite the variance σ_m^2 as

$$\sigma_m^2 = \frac{\sigma^2}{m} \left[1 + 2 \sum_{d=1}^{n-1} \kappa_d \right].$$

The code here shows the evolution of κ_d as a function of d for a series of random numbers. We see that the function κ_d approaches 0 as $d \rightarrow \infty$.

In this case, our data are given by random numbers generated for the uniform distribution with $x \in [0, 1]$. Even with two random numbers being far away, we note that the correlation function is not zero.

Computing the correlation function

This code is best seen with the jupyter-notebook

```
#!/usr/bin/env python
import numpy as np
import matplotlib.mlab as mlab
import matplotlib.pyplot as plt
import random

# initialize the rng with a seed, simple uniform distribution
random.seed()
m = 10000
samplefactor = 1.0/m
x = np.zeros(m)
MeanValue = 0.
VarValue = 0.
for i in range (m):
    value = random.random()
    x[i] = value
    MeanValue += value
    VarValue += value*value

MeanValue *= samplefactor
VarValue *= samplefactor
Variance = VarValue-MeanValue*MeanValue
STDev = np.sqrt(Variance)
print("MeanValue =", MeanValue)
print("Variance =", Variance)
print("Standard deviation =", STDev)
```

Resampling methods: Blocking

The blocking method was made popular by [Flyvbjerg and Pedersen \(1989\)](#) and has become one of the standard ways to estimate the variance $\text{var}(\hat{\theta})$ for exactly one estimator $\hat{\theta}$, namely $\hat{\theta} = \bar{X}$, the mean value.

Assume $n = 2^d$ for some integer $d > 1$ and X_1, X_2, \dots, X_n is a stationary time series to begin with. Moreover, assume that the series is asymptotically uncorrelated. We switch to vector notation by arranging X_1, X_2, \dots, X_n in an n -tuple. Define:

$$\hat{X} = (X_1, X_2, \dots, X_n).$$

Why blocking?

The strength of the blocking method is when the number of observations, n is large. For large n , the complexity of dependent bootstrapping scales poorly, but the blocking method does not, moreover, it becomes more accurate the larger n is.

Blocking Transformations

We now define the blocking transformations. The idea is to take the mean of subsequent pair of elements from \mathbf{X} and form a new vector \mathbf{X}_1 . Continuing in the same way by taking the mean of subsequent pairs of elements of \mathbf{X}_1 we obtain \mathbf{X}_2 , and so on. Define \mathbf{X}_i recursively by:

$$\begin{aligned}(\mathbf{X}_0)_k &\equiv (\mathbf{X})_k \\(\mathbf{X}_{i+1})_k &\equiv \frac{1}{2} \left((\mathbf{X}_i)_{2k-1} + (\mathbf{X}_i)_{2k} \right) \quad \text{for all} \quad 1 \leq i \leq d-1\end{aligned}\tag{1}$$

Blocking transformations

The quantity \mathbf{X}_k is subject to k **blocking transformations**. We now have d vectors $\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_{d-1}$ containing the subsequent averages of observations. It turns out that if the components of \mathbf{X} is a stationary time series, then the components of \mathbf{X}_i is a stationary time series for all $0 \leq i \leq d - 1$

We can then compute the autocovariance (or just covariance), the variance, sample mean, and number of observations for each i . Let $\gamma_i, \sigma_i^2, \bar{X}_i$ denote the covariance, variance and average of the elements of \mathbf{X}_i and let n_i be the number of elements of \mathbf{X}_i . It follows by induction that $n_i = n/2^i$.

Blocking Transformations

Using the definition of the blocking transformation and the distributive property of the covariance, it is clear that since $h = |i - j|$ we can define

$$\begin{aligned}\gamma_{k+1}(h) &= \text{cov}((X_{k+1})_i, (X_{k+1})_j) \\ &= \frac{1}{4} \text{cov}((X_k)_{2i-1} + (X_k)_{2i}, (X_k)_{2j-1} + (X_k)_{2j}) \\ &= \frac{1}{2} \gamma_k(2h) + \frac{1}{2} \gamma_k(2h+1) \quad h = 0\end{aligned}\tag{2}$$

$$= \frac{1}{4} \gamma_k(2h-1) + \frac{1}{2} \gamma_k(2h) + \frac{1}{4} \gamma_k(2h+1) \quad \text{else} \tag{3}$$

The quantity \hat{X} is asymptotically uncorrelated by assumption, \hat{X}_k is also asymptotic uncorrelated. Let's turn our attention to the variance of the sample mean $\text{var}(\bar{X})$.

Blocking Transformations, getting there

We have

$$\text{var}(\bar{X}_k) = \frac{\sigma_k^2}{n_k} + \underbrace{\frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h)}_{\equiv e_k} = \frac{\sigma_k^2}{n_k} + e_k \quad \text{if } \gamma_k(0) = \sigma_k^2. \quad (4)$$

The term e_k is called the **truncation error**:

$$e_k = \frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h). \quad (5)$$

We can show that $\text{var}(\bar{X}_i) = \text{var}(\bar{X}_j)$ for all $0 \leq i \leq d-1$ and $0 \leq j \leq d-1$.

Blocking Transformations, final expressions

We can then wrap up

$$\begin{aligned}n_{j+1}\bar{X}_{j+1} &= \sum_{i=1}^{n_{j+1}} (\hat{X}_{j+1})_i = \frac{1}{2} \sum_{i=1}^{n_j/2} (\hat{X}_j)_{2i-1} + (\hat{X}_j)_{2i} \\&= \frac{1}{2} \left[(\hat{X}_j)_1 + (\hat{X}_j)_2 + \cdots + (\hat{X}_j)_{n_j} \right] = \underbrace{\frac{n_j}{2}}_{=n_{j+1}} \bar{X}_j = n_{j+1}\bar{X}_j.\end{aligned}\tag{6}$$

By repeated use of this equation we get

$\text{var}(\bar{X}_i) = \text{var}(\bar{X}_0) = \text{var}(\bar{X})$ for all $0 \leq i \leq d-1$. This has the consequence that

$$\text{var}(\bar{X}) = \frac{\sigma_k^2}{n_k} + e_k \quad \text{for all} \quad 0 \leq k \leq d-1. \tag{7}$$

More on the blocking method

Flyvbjerg and Petersen demonstrated that the sequence $\{e_k\}_{k=0}^{d-1}$ is decreasing, and conjecture that the term e_k can be made as small as we would like by making k (and hence d) sufficiently large. The sequence is decreasing. It means we can apply blocking transformations until e_k is sufficiently small, and then estimate $\text{var}(\overline{X})$ by $\hat{\sigma}_k^2/n_k$.

For an elegant solution and proof of the blocking method, see the recent article of [Marius Jonsson](#) (former MSc student of the Computational Physics group).

Example code from last week

```
# 2-electron VMC code for 2dim quantum dot with importance sampling
# Using gaussian rng for new positions and Metropolis- Hastings
# Added energy minimization
from math import exp, sqrt
from random import random, seed, normalvariate
import numpy as np
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from matplotlib import cm
from matplotlib.ticker import LinearLocator, FormatStrFormatter
from scipy.optimize import minimize
import sys
import os

# Where to save data files
PROJECT_ROOT_DIR = "Results"
DATA_ID = "Results/EnergyMin"

if not os.path.exists(PROJECT_ROOT_DIR):
    os.mkdir(PROJECT_ROOT_DIR)

if not os.path.exists(DATA_ID):
    os.makedirs(DATA_ID)

def data_path(dat_id):
    return os.path.join(DATA_ID, dat_id)

outfile = open(data_path("Energies.dat"), 'w')
```

Resampling analysis

The next step is then to use the above data sets and perform a resampling analysis using the blocking method The blocking code, based on the article of [Marius Jonsson](#) is given here

```
# Common imports
```

```
import os
```

```
# Where to save the figures and data files
```

```
DATA_ID = "Results/EnergyMin"
```

```
def data_path(dat_id):
```

```
    return os.path.join(DATA_ID, dat_id)
```

```
infile = open(data_path("Energies.dat"), 'r')
```

```
from numpy import log2, zeros, mean, var, sum, loadtxt, arange, array,
```

```
from numpy.linalg import inv
```

```
def block(x):
```

```
    # preliminaries
```

```
    n = len(x)
```

```
    d = int(log2(n))
```

```
    s, gamma = zeros(d), zeros(d)
```

```
    mu = mean(x)
```

```
    # estimate the auto-covariance and variances
```

```
    # for each blocking transformation
```

```
    for i in arange(0, d):
```