

# Week 10, March 4-8: Resampling Techniques, Bootstrap and Blocking

Morten Hjorth-Jensen Email [morten.hjorth-jensen@fys.uio.no](mailto:morten.hjorth-jensen@fys.uio.no)<sup>1,2</sup>

<sup>1</sup>Department of Physics and Center for Computing in Science Education, University of Oslo, Oslo, Norway

<sup>2</sup>Department of Physics and Astronomy and Facility for Rare Ion Beams, Michigan State University, East Lansing, Michigan, USA

March 4-8

## Overview of week 10, March 4-8

### Topics.

- Reminder from last week on optimization methods
- Resampling Techniques and statistics: Bootstrap and Blocking
- [Video of lecture to be added](#)
- [Handwritten notes](#)

### Teaching Material, videos and written material.

- Overview video on the [Bootstrap method](#)
- [Marius Johnson's Master thesis on the Blocking Method](#)

## Resampling methods

Resampling methods are an indispensable tool in modern statistics. They involve repeatedly drawing samples from a training set and refitting a model of interest on each sample in order to obtain additional information about the fitted model. For example, in order to estimate the variability of a linear regression fit, we can repeatedly draw different samples from the training data, fit a linear regression to each new sample, and then examine the extent to which the resulting fits differ. Such an approach may allow us to obtain information that would not be available from fitting the model only once using the original training sample.

## Resampling approaches can be computationally expensive

Resampling approaches can be computationally expensive, because they involve fitting the same statistical method multiple times using different subsets of the training data. However, due to recent advances in computing power, the computational requirements of resampling methods generally are not prohibitive. In this chapter, we discuss two of the most commonly used resampling methods, cross-validation and the bootstrap. Both methods are important tools in the practical application of many statistical learning procedures. For example, cross-validation can be used to estimate the test error associated with a given statistical learning method in order to evaluate its performance, or to select the appropriate level of flexibility. The process of evaluating a model's performance is known as model assessment, whereas the process of selecting the proper level of flexibility for a model is known as model selection. The bootstrap is widely used.

## 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:
  - Statistical errors
  - 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.

## Statistics

The *probability distribution function (PDF)* is a function  $p(x)$  on the domain which, in the discrete case, gives us the probability or relative frequency with which these values of  $X$  occur:

$$p(x) = \text{prob}(X = x)$$

In the continuous case, the PDF does not directly depict the actual probability. Instead we define the probability for the stochastic variable to assume any value on an infinitesimal interval around  $x$  to be  $p(x)dx$ . The continuous function  $p(x)$  then gives us the *density* of the probability rather than the probability itself. The probability for a stochastic variable to assume any value on a non-infinitesimal interval  $[a, b]$  is then just the integral:

$$\text{prob}(a \leq X \leq b) = \int_a^b p(x)dx$$

Qualitatively speaking, a stochastic variable represents the values of numbers chosen as if by chance from some specified PDF so that the selection of a large set of these numbers reproduces this PDF.

## Statistics, moments

A particularly useful class of special expectation values are the *moments*. The  $n$ -th moment of the PDF  $p$  is defined as follows:

$$\langle x^n \rangle \equiv \int x^n p(x) dx$$

The zero-th moment  $\langle 1 \rangle$  is just the normalization condition of  $p$ . The first moment,  $\langle x \rangle$ , is called the *mean* of  $p$  and often denoted by the letter  $\mu$ :

$$\langle x \rangle = \mu \equiv \int x p(x) dx$$

## Statistics, central moments

A special version of the moments is the set of *central moments*, the  $n$ -th central moment defined as:

$$\langle (x - \langle x \rangle)^n \rangle \equiv \int (x - \langle x \rangle)^n p(x) dx$$

The zero-th and first central moments are both trivial, equal 1 and 0, respectively. But the second central moment, known as the *variance* of  $p$ , is of particular

interest. For the stochastic variable  $X$ , the variance is denoted as  $\sigma_X^2$  or  $\text{var}(X)$ :

$$\sigma_X^2 = \text{var}(X) = \langle (x - \langle x \rangle)^2 \rangle = \int (x - \langle x \rangle)^2 p(x) dx \quad (1)$$

$$= \int (x^2 - 2x\langle x \rangle + \langle x \rangle^2) p(x) dx \quad (2)$$

$$= \langle x^2 \rangle - 2\langle x \rangle \langle x \rangle + \langle x \rangle^2 \quad (3)$$

$$= \langle x^2 \rangle - \langle x \rangle^2 \quad (4)$$

The square root of the variance,  $\sigma = \sqrt{\langle (x - \langle x \rangle)^2 \rangle}$  is called the *standard deviation* of  $p$ . It is clearly just the RMS (root-mean-square) value of the deviation of the PDF from its mean value, interpreted qualitatively as the *spread* of  $p$  around its mean.

## Statistics, covariance

Another important quantity is the so called covariance, a variant of the above defined variance. Consider again the set  $\{X_i\}$  of  $n$  stochastic variables (not necessarily uncorrelated) with the multivariate PDF  $P(x_1, \dots, x_n)$ . The *covariance* of two of the stochastic variables,  $X_i$  and  $X_j$ , is defined as follows:

$$\begin{aligned} \text{cov}(X_i, X_j) &\equiv \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \\ &= \int \dots \int (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) P(x_1, \dots, x_n) dx_1 \dots dx_n \end{aligned} \quad (5)$$

with

$$\langle x_i \rangle = \int \dots \int x_i P(x_1, \dots, x_n) dx_1 \dots dx_n$$

## Statistics, more covariance

If we consider the above covariance as a matrix  $C_{ij} = \text{cov}(X_i, X_j)$ , then the diagonal elements are just the familiar variances,  $C_{ii} = \text{cov}(X_i, X_i) = \text{var}(X_i)$ . It turns out that all the off-diagonal elements are zero if the stochastic variables are uncorrelated. This is easy to show, keeping in mind the linearity of the expectation value. Consider the stochastic variables  $X_i$  and  $X_j$ , ( $i \neq j$ ):

$$\text{cov}(X_i, X_j) = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \quad (6)$$

$$= \langle x_i x_j - x_i \langle x_j \rangle - \langle x_i \rangle x_j + \langle x_i \rangle \langle x_j \rangle \rangle \quad (7)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle - \langle \langle x_i \rangle x_j \rangle + \langle \langle x_i \rangle \langle x_j \rangle \rangle \quad (8)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle - \langle x_i \rangle \langle x_j \rangle + \langle x_i \rangle \langle x_j \rangle \quad (9)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle \quad (10)$$

## Statistics, independent variables

If  $X_i$  and  $X_j$  are independent, we get  $\langle x_i x_j \rangle = \langle x_i \rangle \langle x_j \rangle$ , resulting in  $\text{cov}(X_i, X_j) = 0$  ( $i \neq j$ ).

Also useful for us is the covariance of linear combinations of stochastic variables. Let  $\{X_i\}$  and  $\{Y_i\}$  be two sets of stochastic variables. Let also  $\{a_i\}$  and  $\{b_i\}$  be two sets of scalars. Consider the linear combination:

$$U = \sum_i a_i X_i \quad V = \sum_j b_j Y_j$$

By the linearity of the expectation value

$$\text{cov}(U, V) = \sum_{i,j} a_i b_j \text{cov}(X_i, Y_j)$$

### Statistics, more variance

Now, since the variance is just  $\text{var}(X_i) = \text{cov}(X_i, X_i)$ , we get the variance of the linear combination  $U = \sum_i a_i X_i$ :

$$\text{var}(U) = \sum_{i,j} a_i a_j \text{cov}(X_i, X_j) \quad (11)$$

And in the special case when the stochastic variables are uncorrelated, the off-diagonal elements of the covariance are as we know zero, resulting in:

$$\begin{aligned} \text{var}(U) &= \sum_i a_i^2 \text{cov}(X_i, X_i) = \sum_i a_i^2 \text{var}(X_i) \\ \text{var}\left(\sum_i a_i X_i\right) &= \sum_i a_i^2 \text{var}(X_i) \end{aligned}$$

which will become very useful in our study of the error in the mean value of a set of measurements.

### Statistics and stochastic processes

A *stochastic process* is a process that produces sequentially a chain of values:

$$\{x_1, x_2, \dots, x_k, \dots\}.$$

We will call these values our *measurements* and the entire set as our measured *sample*. The action of measuring all the elements of a sample we will call a stochastic *experiment* since, operationally, they are often associated with results of empirical observation of some physical or mathematical phenomena; precisely an experiment. We assume that these values are distributed according to some PDF  $p_X(x)$ , where  $X$  is just the formal symbol for the stochastic variable whose PDF is  $p_X(x)$ . Instead of trying to determine the full distribution  $p$  we are often only interested in finding the few lowest moments, like the mean  $\mu_X$  and the variance  $\sigma_X$ .

## Statistics and sample variables

In practical situations a sample is always of finite size. Let that size be  $n$ . The expectation value of a sample, the *sample mean*, is then defined as follows:

$$\bar{x}_n \equiv \frac{1}{n} \sum_{k=1}^n x_k$$

The *sample variance* is:

$$\text{var}(x) \equiv \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n)^2$$

its square root being the *standard deviation of the sample*. The *sample covariance* is:

$$\text{cov}(x) \equiv \frac{1}{n} \sum_{kl} (x_k - \bar{x}_n)(x_l - \bar{x}_n)$$

## Statistics, sample variance and covariance

Note that the sample variance is the sample covariance without the cross terms. In a similar manner as the covariance in Eq. (5) is a measure of the correlation between two stochastic variables, the above defined sample covariance is a measure of the sequential correlation between succeeding measurements of a sample.

These quantities, being known experimental values, differ significantly from and must not be confused with the similarly named quantities for stochastic variables, mean  $\mu_X$ , variance  $\text{var}(X)$  and covariance  $\text{cov}(X, Y)$ .

## Statistics, law of large numbers

The law of large numbers states that as the size of our sample grows to infinity, the sample mean approaches the true mean  $\mu_X$  of the chosen PDF:

$$\lim_{n \rightarrow \infty} \bar{x}_n = \mu_X$$

The sample mean  $\bar{x}_n$  works therefore as an estimate of the true mean  $\mu_X$ .

What we need to find out is how good an approximation  $\bar{x}_n$  is to  $\mu_X$ . In any stochastic measurement, an estimated mean is of no use to us without a measure of its error. A quantity that tells us how well we can reproduce it in another experiment. We are therefore interested in the PDF of the sample mean itself. Its standard deviation will be a measure of the spread of sample means, and we will simply call it the *error* of the sample mean, or just sample error, and denote it by  $\text{err}_X$ . In practice, we will only be able to produce an *estimate* of the sample error since the exact value would require the knowledge of the true PDFs behind, which we usually do not have.

## Statistics, more on sample error

Let us first take a look at what happens to the sample error as the size of the sample grows. In a sample, each of the measurements  $x_i$  can be associated with its own stochastic variable  $X_i$ . The stochastic variable  $\bar{X}_n$  for the sample mean  $\bar{x}_n$  is then just a linear combination, already familiar to us:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

All the coefficients are just equal  $1/n$ . The PDF of  $\bar{X}_n$ , denoted by  $p_{\bar{X}_n}(x)$  is the desired PDF of the sample means.

## Statistics

The probability density of obtaining a sample mean  $\bar{x}_n$  is the product of probabilities of obtaining arbitrary values  $x_1, x_2, \dots, x_n$  with the constraint that the mean of the set  $\{x_i\}$  is  $\bar{x}_n$ :

$$p_{\bar{X}_n}(x) = \int p_X(x_1) \cdots \int p_X(x_n) \delta\left(x - \frac{x_1 + x_2 + \cdots + x_n}{n}\right) dx_n \cdots dx_1$$

And in particular we are interested in its variance  $\text{var}(\bar{X}_n)$ .

## Statistics, central limit theorem

It is generally not possible to express  $p_{\bar{X}_n}(x)$  in a closed form given an arbitrary PDF  $p_X$  and a number  $n$ . But for the limit  $n \rightarrow \infty$  it is possible to make an approximation. The very important result is called *the central limit theorem*. It tells us that as  $n$  goes to infinity,  $p_{\bar{X}_n}(x)$  approaches a Gaussian distribution whose mean and variance equal the true mean and variance,  $\mu_X$  and  $\sigma_X^2$ , respectively:

$$\lim_{n \rightarrow \infty} p_{\bar{X}_n}(x) = \left( \frac{n}{2\pi \text{var}(X)} \right)^{1/2} e^{-\frac{n(x - \bar{x}_n)^2}{2\text{var}(X)}} \quad (12)$$

## Statistics, more technicalities

The desired variance  $\text{var}(\bar{X}_n)$ , i.e. the sample error squared  $\text{err}_X^2$ , is given by:

$$\text{err}_X^2 = \text{var}(\bar{X}_n) = \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) \quad (13)$$

We see now that in order to calculate the exact error of the sample with the above expression, we would need the true means  $\mu_{X_i}$  of the stochastic variables  $X_i$ . To calculate these requires that we know the true multivariate PDF of all the  $X_i$ . But this PDF is unknown to us, we have only got the measurements of

one sample. The best we can do is to let the sample itself be an estimate of the PDF of each of the  $X_i$ , estimating all properties of  $X_i$  through the measurements of the sample.

## Statistics

Our estimate of  $\mu_{X_i}$  is then the sample mean  $\bar{x}$  itself, in accordance with the central limit theorem:

$$\mu_{X_i} = \langle x_i \rangle \approx \frac{1}{n} \sum_{k=1}^n x_k = \bar{x}$$

Using  $\bar{x}$  in place of  $\mu_{X_i}$  we can give an *estimate* of the covariance in Eq. (13)

$$\text{cov}(X_i, X_j) = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \approx \langle (x_i - \bar{x})(x_j - \bar{x}) \rangle,$$

resulting in

$$\frac{1}{n} \sum_l^n \left( \frac{1}{n} \sum_k^n (x_k - \bar{x}_n)(x_l - \bar{x}_n) \right) = \frac{1}{n} \frac{1}{n} \sum_{kl} (x_k - \bar{x}_n)(x_l - \bar{x}_n) = \frac{1}{n} \text{cov}(x)$$

## Statistics and sample variance

By the same procedure we can use the sample variance as an estimate of the variance of any of the stochastic variables  $X_i$

$$\text{var}(X_i) = \langle x_i - \langle x_i \rangle \rangle \approx \langle x_i - \bar{x}_n \rangle,$$

which is approximated as

$$\text{var}(X_i) \approx \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n) = \text{var}(x) \quad (14)$$

Now we can calculate an estimate of the error  $\text{err}_X$  of the sample mean  $\bar{x}_n$ :

$$\begin{aligned} \text{err}_X^2 &= \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) \\ &\approx \frac{1}{n^2} \sum_{ij} \frac{1}{n} \text{cov}(x) = \frac{1}{n^2} n^2 \frac{1}{n} \text{cov}(x) \\ &= \frac{1}{n} \text{cov}(x) \end{aligned} \quad (15)$$

which is nothing but the sample covariance divided by the number of measurements in the sample.



## Statistics, uncorrelated results

In the special case that the measurements of the sample are uncorrelated (equivalently the stochastic variables  $X_i$  are uncorrelated) we have that the off-diagonal elements of the covariance are zero. This gives the following estimate of the sample error:

$$\text{err}_X^2 = \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) = \frac{1}{n^2} \sum_i \text{var}(X_i),$$

resulting in

$$\text{err}_X^2 \approx \frac{1}{n^2} \sum_i \text{var}(x) = \frac{1}{n} \text{var}(x) \quad (16)$$

where in the second step we have used Eq. (14). The error of the sample is then just its standard deviation divided by the square root of the number of measurements the sample contains. This is a very useful formula which is easy to compute. It acts as a first approximation to the error, but in numerical experiments, we cannot overlook the always present correlations.

## Statistics, computations

For computational purposes one usually splits up the estimate of  $\text{err}_X^2$ , given by Eq. (15), into two parts

$$\text{err}_X^2 = \frac{1}{n} \text{var}(x) + \frac{1}{n} (\text{cov}(x) - \text{var}(x)),$$

which equals

$$\frac{1}{n^2} \sum_{k=1}^n (x_k - \bar{x}_n)^2 + \frac{2}{n^2} \sum_{k < l} (x_k - \bar{x}_n)(x_l - \bar{x}_n) \quad (17)$$

The first term is the same as the error in the uncorrelated case, Eq. (16). This means that the second term accounts for the error correction due to correlation between the measurements. For uncorrelated measurements this second term is zero.

## Statistics, more on computations of errors

Computationally the uncorrelated first term is much easier to treat efficiently than the second.

$$\text{var}(x) = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n)^2 = \left( \frac{1}{n} \sum_{k=1}^n x_k^2 \right) - \bar{x}_n^2$$

We just accumulate separately the values  $x^2$  and  $x$  for every measurement  $x$  we receive. The correlation term, though, has to be calculated at the end of the experiment since we need all the measurements to calculate the cross terms. Therefore, all measurements have to be stored throughout the experiment.

## Statistics, wrapping up 1

Let us analyze the problem by splitting up the correlation term into partial sums of the form:

$$f_d = \frac{1}{n-d} \sum_{k=1}^{n-d} (x_k - \bar{x}_n)(x_{k+d} - \bar{x}_n)$$

The correlation term of the error can now be rewritten in terms of  $f_d$

$$\frac{2}{n} \sum_{k < l} (x_k - \bar{x}_n)(x_l - \bar{x}_n) = 2 \sum_{d=1}^{n-1} f_d$$

The value of  $f_d$  reflects the correlation between measurements separated by the distance  $d$  in the sample samples. Notice that for  $d = 0$ ,  $f$  is just the sample variance,  $\text{var}(x)$ . If we divide  $f_d$  by  $\text{var}(x)$ , we arrive at the so called *autocorrelation function*

$$\kappa_d = \frac{f_d}{\text{var}(x)}$$

which gives us a useful measure of pairwise correlations starting always at 1 for  $d = 0$ .

## Statistics, final expression

The sample error (see eq. (17)) can now be written in terms of the autocorrelation function:

$$\begin{aligned} \text{err}_X^2 &= \frac{1}{n} \text{var}(x) + \frac{2}{n} \cdot \text{var}(x) \sum_{d=1}^{n-1} \frac{f_d}{\text{var}(x)} \\ &= \left( 1 + 2 \sum_{d=1}^{n-1} \kappa_d \right) \frac{1}{n} \text{var}(x) \\ &= \frac{\tau}{n} \cdot \text{var}(x) \end{aligned} \tag{18}$$

and we see that  $\text{err}_X$  can be expressed in terms the uncorrelated sample variance times a correction factor  $\tau$  which accounts for the correlation between measurements. We call this correction factor the *autocorrelation time*:

$$\tau = 1 + 2 \sum_{d=1}^{n-1} \kappa_d \tag{19}$$

## Statistics, effective number of correlations

For a correlation free experiment,  $\tau$  equals 1. From the point of view of eq. (18) we can interpret a sequential correlation as an effective reduction of the number of measurements by a factor  $\tau$ . The effective number of measurements becomes:

$$n_{\text{eff}} = \frac{n}{\tau}$$

To neglect the autocorrelation time  $\tau$  will always cause our simple uncorrelated estimate of  $\text{err}_X^2 \approx \text{var}(x)/n$  to be less than the true sample error. The estimate of the error will be too *good*. On the other hand, the calculation of the full autocorrelation time poses an efficiency problem if the set of measurements is very large.

## Can we understand this? Time Auto-correlation Function

The so-called time-displacement autocorrelation  $\phi(t)$  for a quantity  $\mathbf{M}$  is given by

$$\phi(t) = \int dt' [\mathbf{M}(t') - \langle \mathbf{M} \rangle] [\mathbf{M}(t' + t) - \langle \mathbf{M} \rangle],$$

which can be rewritten as

$$\phi(t) = \int dt' [\mathbf{M}(t')\mathbf{M}(t' + t) - \langle \mathbf{M} \rangle^2],$$

where  $\langle \mathbf{M} \rangle$  is the average value and  $\mathbf{M}(t)$  its instantaneous value. We can discretize this function as follows, where we used our set of computed values  $\mathbf{M}(t)$  for a set of discretized times (our Monte Carlo cycles corresponding to moving all electrons?)

$$\phi(t) = \frac{1}{t_{\max} - t} \sum_{t'=0}^{t_{\max}-t} \mathbf{M}(t')\mathbf{M}(t'+t) - \frac{1}{t_{\max} - t} \sum_{t'=0}^{t_{\max}-t} \mathbf{M}(t') \times \frac{1}{t_{\max} - t} \sum_{t'=0}^{t_{\max}-t} \mathbf{M}(t'+t).$$

## Time Auto-correlation Function

One should be careful with times close to  $t_{\max}$ , the upper limit of the sums becomes small and we end up integrating over a rather small time interval. This means that the statistical error in  $\phi(t)$  due to the random nature of the fluctuations in  $\mathbf{M}(t)$  can become large.

One should therefore choose  $t \ll t_{\max}$ .

Note that the variable  $\mathbf{M}$  can be any expectation values of interest.

The time-correlation function gives a measure of the correlation between the various values of the variable at a time  $t'$  and a time  $t' + t$ . If we multiply the values of  $\mathbf{M}$  at these two different times, we will get a positive contribution if they are fluctuating in the same direction, or a negative value if they fluctuate in the opposite direction. If we then integrate over time, or use the discretized version of, the time correlation function  $\phi(t)$  should take a non-zero value if the fluctuations are correlated, else it should gradually go to zero. For times a long way apart the different values of  $\mathbf{M}$  are most likely uncorrelated and  $\phi(t)$  should be zero.

## Time Auto-correlation Function

We can derive the correlation time by observing that our Metropolis algorithm is based on a random walk in the space of all possible spin configurations. Our

probability distribution function  $\hat{\mathbf{w}}(t)$  after a given number of time steps  $t$  could be written as

$$\hat{\mathbf{w}}(t) = \hat{\mathbf{W}}^t \hat{\mathbf{w}}(0),$$

with  $\hat{\mathbf{w}}(0)$  the distribution at  $t = 0$  and  $\hat{\mathbf{W}}$  representing the transition probability matrix. We can always expand  $\hat{\mathbf{w}}(0)$  in terms of the right eigenvectors of  $\hat{\mathbf{W}}$  as

$$\hat{\mathbf{w}}(0) = \sum_i \alpha_i \hat{\mathbf{v}}_i,$$

resulting in

$$\hat{\mathbf{w}}(t) = \hat{\mathbf{W}}^t \hat{\mathbf{w}}(0) = \hat{\mathbf{W}}^t \sum_i \alpha_i \hat{\mathbf{v}}_i = \sum_i \lambda_i^t \alpha_i \hat{\mathbf{v}}_i,$$

with  $\lambda_i$  the  $i^{\text{th}}$  eigenvalue corresponding to the eigenvector  $\hat{\mathbf{v}}_i$ .

## Time Auto-correlation Function

If we assume that  $\lambda_0$  is the largest eigenvalue we see that in the limit  $t \rightarrow \infty$ ,  $\hat{\mathbf{w}}(t)$  becomes proportional to the corresponding eigenvector  $\hat{\mathbf{v}}_0$ . This is our steady state or final distribution.

We can relate this property to an observable like the mean energy. With the probability  $\hat{\mathbf{w}}(t)$  (which in our case is the squared trial wave function) we can write the expectation values as

$$\langle \mathbf{M}(t) \rangle = \sum_{\mu} \hat{\mathbf{w}}(t)_{\mu} \mathbf{M}_{\mu},$$

or as the scalar of a vector product

$$\langle \mathbf{M}(t) \rangle = \hat{\mathbf{w}}(t) \mathbf{m},$$

with  $\mathbf{m}$  being the vector whose elements are the values of  $\mathbf{M}_{\mu}$  in its various microstates  $\mu$ .

## Time Auto-correlation Function

We rewrite this relation as

$$\langle \mathbf{M}(t) \rangle = \hat{\mathbf{w}}(t) \mathbf{m} = \sum_i \lambda_i^t \alpha_i \hat{\mathbf{v}}_i \mathbf{m}_i.$$

If we define  $m_i = \hat{\mathbf{v}}_i \mathbf{m}_i$  as the expectation value of  $\mathbf{M}$  in the  $i^{\text{th}}$  eigenstate we can rewrite the last equation as

$$\langle \mathbf{M}(t) \rangle = \sum_i \lambda_i^t \alpha_i m_i.$$

Since we have that in the limit  $t \rightarrow \infty$  the mean value is dominated by the the largest eigenvalue  $\lambda_0$ , we can rewrite the last equation as

$$\langle \mathbf{M}(t) \rangle = \langle \mathbf{M}(\infty) \rangle + \sum_{i \neq 0} \lambda_i^t \alpha_i m_i.$$

We define the quantity

$$\tau_i = -\frac{1}{\log \lambda_i},$$

and rewrite the last expectation value as

$$\langle \mathbf{M}(t) \rangle = \langle \mathbf{M}(\infty) \rangle + \sum_{i \neq 0} \alpha_i m_i e^{-t/\tau_i}.$$

## Time Auto-correlation Function

The quantities  $\tau_i$  are the correlation times for the system. They control also the auto-correlation function discussed above. The longest correlation time is obviously given by the second largest eigenvalue  $\tau_1$ , which normally defines the correlation time discussed above. For large times, this is the only correlation time that survives. If higher eigenvalues of the transition matrix are well separated from  $\lambda_1$  and we simulate long enough,  $\tau_1$  may well define the correlation time. In other cases we may not be able to extract a reliable result for  $\tau_1$ . Coming back to the time correlation function  $\phi(t)$  we can present a more general definition in terms of the mean magnetizations  $\langle \mathbf{M}(t) \rangle$ . Recalling that the mean value is equal to  $\langle \mathbf{M}(\infty) \rangle$  we arrive at the expectation values

$$\phi(t) = \langle \mathbf{M}(0) - \mathbf{M}(\infty) \rangle \langle \mathbf{M}(t) - \mathbf{M}(\infty) \rangle,$$

resulting in

$$\phi(t) = \sum_{i,j \neq 0} m_i \alpha_i m_j \alpha_j e^{-t/\tau_i},$$

which is appropriate for all times.

## Correlation Time

If the correlation function decays exponentially

$$\phi(t) \sim \exp(-t/\tau)$$

then the exponential correlation time can be computed as the average

$$\tau_{\text{exp}} = -\left\langle \frac{t}{\log \left| \frac{\phi(t)}{\phi(0)} \right|} \right\rangle.$$

If the decay is exponential, then

$$\int_0^\infty dt \phi(t) = \int_0^\infty dt \phi(0) \exp(-t/\tau) = \tau \phi(0),$$

which suggests another measure of correlation

$$\tau_{\text{int}} = \sum_k \frac{\phi(k)}{\phi(0)},$$

called the integrated correlation time.

## Resampling methods: Jackknife and Bootstrap

Two famous resampling methods are the **independent bootstrap** and **the jackknife**.

The jackknife is a special case of the independent bootstrap. Still, the jackknife was made popular prior to the independent bootstrap. And as the popularity of the independent bootstrap soared, new variants, such as **the dependent bootstrap**.

The Jackknife and independent bootstrap work for independent, identically distributed random variables. If these conditions are not satisfied, the methods will fail. Yet, it should be said that if the data are independent, identically distributed, and we only want to estimate the variance of  $\bar{X}$  (which often is the case), then there is no need for bootstrapping.

## Resampling methods: Jackknife

The Jackknife works by making many replicas of the estimator  $\hat{\theta}$ . The jackknife is a resampling method, we explained that this happens by scrambling the data in some way. When using the jackknife, this is done by systematically leaving out one observation from the vector of observed values  $\hat{x} = (x_1, x_2, \dots, x_n)$ . Let  $\hat{x}_i$  denote the vector

$$\hat{x}_i = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

which equals the vector  $\hat{x}$  with the exception that observation number  $i$  is left out. Using this notation, define  $\hat{\theta}_i$  to be the estimator  $\hat{\theta}$  computed using  $\vec{X}_i$ .

## Resampling methods: Jackknife estimator

To get an estimate for the bias and standard error of  $\hat{\theta}$ , use the following estimators for each component of  $\hat{\theta}$

$$\widehat{\text{Bias}}(\hat{\theta}, \theta) = (n-1) \left( -\hat{\theta} + \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i \right) \quad \text{and} \quad \hat{\sigma}_{\hat{\theta}}^2 = \frac{n-1}{n} \sum_{i=1}^n \left( \hat{\theta}_i - \frac{1}{n} \sum_{j=1}^n \hat{\theta}_j \right)^2.$$

## Jackknife code example

```
from numpy import *
from numpy.random import randint, randn
```

```

from time import time

def jackknife(data, stat):
    n = len(data); t = zeros(n); inds = arange(n); t0 = time()
    ## 'jackknifing' by leaving out an observation for each i
    for i in range(n):
        t[i] = stat(delete(data,i) )

    # analysis
    print("Runtime: %g sec" % (time()-t0)); print("Jackknife Statistics :")
    print("original      bias      std. error")
    print("%8g %14g %15g" % (stat(data), (n-1)*mean(t)-stat(data), ((n-1)*var(t))**.5))

    return t

# Returns mean of data samples
def stat(data):
    return mean(data)

mu, sigma = 100, 15
datapoints = 10000
x = mu + sigma*random.randn(datapoints)
# jackknife returns the data sample
t = jackknife(x, stat)

```

## Resampling methods: Bootstrap

Bootstrapping is a nonparametric approach to statistical inference that substitutes computation for more traditional distributional assumptions and asymptotic results. Bootstrapping offers a number of advantages:

1. The bootstrap is quite general, although there are some cases in which it fails.
2. Because it does not require distributional assumptions (such as normally distributed errors), the bootstrap can provide more accurate inferences when the data are not well behaved or when the sample size is small.
3. It is possible to apply the bootstrap to statistics with sampling distributions that are difficult to derive, even asymptotically.
4. It is relatively simple to apply the bootstrap to complex data-collection plans (such as stratified and clustered samples).

## Resampling methods: Bootstrap background

Since  $\hat{\theta} = \hat{\theta}(\hat{X})$  is a function of random variables,  $\hat{\theta}$  itself must be a random variable. Thus it has a pdf, call this function  $p(\hat{t})$ . The aim of the bootstrap is to estimate  $p(\hat{t})$  by the relative frequency of  $\hat{\theta}$ . You can think of this as using

a histogram in the place of  $p(\hat{t})$ . If the relative frequency closely resembles  $p(\hat{t})$ , then using numerics, it is straight forward to estimate all the interesting parameters of  $p(\hat{t})$  using point estimators.

## Resampling methods: More Bootstrap background

In the case that  $\hat{\theta}$  has more than one component, and the components are independent, we use the same estimator on each component separately. If the probability density function of  $X_i$ ,  $p(x)$ , had been known, then it would have been straight forward to do this by:

1. Drawing lots of numbers from  $p(x)$ , suppose we call one such set of numbers  $(X_1^*, X_2^*, \dots, X_n^*)$ .
2. Then using these numbers, we could compute a replica of  $\hat{\theta}$  called  $\hat{\theta}^*$ .

By repeated use of (1) and (2), many estimates of  $\hat{\theta}$  could have been obtained. The idea is to use the relative frequency of  $\hat{\theta}^*$  (think of a histogram) as an estimate of  $p(\hat{t})$ .

## Resampling methods: Bootstrap approach

But unless there is enough information available about the process that generated  $X_1, X_2, \dots, X_n$ ,  $p(x)$  is in general unknown. Therefore, [Efron in 1979](#) asked the question: What if we replace  $p(x)$  by the relative frequency of the observation  $X_i$ ; if we draw observations in accordance with the relative frequency of the observations, will we obtain the same result in some asymptotic sense? The answer is yes.

Instead of generating the histogram for the relative frequency of the observation  $X_i$ , just draw the values  $(X_1^*, X_2^*, \dots, X_n^*)$  with replacement from the vector  $\hat{X}$ .

## Resampling methods: Bootstrap steps

The independent bootstrap works like this:

1. Draw with replacement  $n$  numbers for the observed variables  $\hat{x} = (x_1, x_2, \dots, x_n)$ .
2. Define a vector  $\hat{x}^*$  containing the values which were drawn from  $\hat{x}$ .
3. Using the vector  $\hat{x}^*$  compute  $\hat{\theta}^*$  by evaluating  $\hat{\theta}$  under the observations  $\hat{x}^*$ .
4. Repeat this process  $k$  times.

When you are done, you can draw a histogram of the relative frequency of  $\hat{\theta}^*$ . This is your estimate of the probability distribution  $p(t)$ . Using this probability distribution you can estimate any statistics thereof. In principle you never draw the histogram of the relative frequency of  $\hat{\theta}^*$ . Instead you use the estimators corresponding to the statistic of interest. For example, if you are interested in estimating the variance of  $\hat{\theta}$ , apply the estimator  $\hat{\sigma}^2$  to the values  $\hat{\theta}^*$ .



## Code example for the Bootstrap method

The following code starts with a Gaussian distribution with mean value  $\mu = 100$  and variance  $\sigma = 15$ . We use this to generate the data used in the bootstrap analysis. The bootstrap analysis returns a data set after a given number of bootstrap operations (as many as we have data points). This data set consists of estimated mean values for each bootstrap operation. The histogram generated by the bootstrap method shows that the distribution for these mean values is also a Gaussian, centered around the mean value  $\mu = 100$  but with standard deviation  $\sigma/\sqrt{n}$ , where  $n$  is the number of bootstrap samples (in this case the same as the number of original data points). The value of the standard deviation is what we expect from the central limit theorem.

```
%matplotlib inline

from numpy import *
from numpy.random import randint, randn
from time import time
from scipy.stats import norm
import matplotlib.pyplot as plt

# Returns mean of bootstrap samples
def stat(data):
    return mean(data)

# Bootstrap algorithm
def bootstrap(data, statistic, R):
    t = zeros(R); n = len(data); inds = arange(n); t0 = time()

    # non-parametric bootstrap
    for i in range(R):
        t[i] = statistic(data[randint(0,n,n)])

    # analysis
    print("Runtime: %g sec" % (time()-t0)); print("Bootstrap Statistics :")
    print("original      bias      std. error")
    print("%8g %8g %14g %15g" % (statistic(data), std(data), \
                                mean(t), \
                                std(t)))

    return t

mu, sigma = 100, 15
datapoints = 10000
x = mu + sigma*random.randn(datapoints)
# bootstrap returns the data sample
# the histogram of the bootstrapped data
t = bootstrap(x, stat, datapoints)
# the histogram of the bootstrapped data
n, binsboot, patches = plt.hist(t, bins=50, density='true', histtype='bar', color='red', alpha=0.75)

# add a 'best fit' line
y = norm.pdf(binsboot, mean(t), std(t))
lt = plt.plot(binsboot, y, 'r--', linewidth=1)
plt.xlabel('Smarts')
plt.ylabel('Probability')
plt.axis([99.5, 100.6, 0, 3.0])
```

```
plt.grid(True)

plt.show()
```

## Resampling methods: Blocking

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

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 time series is asymptotically uncorrelated. We switch to vector notation by arranging  $X_1, X_2, \dots, X_n$  in an  $n$ -tuple. Define:

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

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

$$\begin{aligned} (\vec{X}_0)_k &\equiv (\vec{X})_k \\ (\vec{X}_{i+1})_k &\equiv \frac{1}{2} \left( (\vec{X}_i)_{2k-1} + (\vec{X}_i)_{2k} \right) \quad \text{for all } 1 \leq i \leq d-1 \end{aligned} \quad (20)$$

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

We can then compute the autocovariance, the variance, sample mean, and number of observations for each  $i$ . Let  $\gamma_i, \sigma_i^2, \bar{X}_i$  denote the autocovariance, variance and average of the elements of  $\vec{X}_i$  and let  $n_i$  be the number of elements of  $\vec{X}_i$ . 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

$$\begin{aligned}\gamma_{k+1}(h) &= \text{cov}((X_{k+1})_i, (X_{k+1})_j) \\ &= \frac{1}{4} \text{cov}((X_k)_{2i-1} + (X_k)_{2i}, (X_k)_{2j-1} + (X_k)_{2j}) \\ &= \frac{1}{2} \gamma_k(2h) + \frac{1}{2} \gamma_k(2h+1) \quad h = 0 \\ &= \frac{1}{4} \gamma_k(2h-1) + \frac{1}{2} \gamma_k(2h) + \frac{1}{4} \gamma_k(2h+1) \quad \text{else}\end{aligned}\tag{21}$$

$$\tag{22}$$

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

## Blocking Transformations, getting there

We have

$$V(\bar{X}_k) = \frac{\sigma_k^2}{n_k} + \underbrace{\frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h)}_{\equiv e_k} = \frac{\sigma_k^2}{n_k} + e_k \quad \text{if } \gamma_k(0) = \sigma_k^2. \tag{23}$$

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

$$e_k = \frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h). \tag{24}$$

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

## Blocking Transformations, final expressions

We can then wrap up

$$\begin{aligned}n_{j+1} \bar{X}_{j+1} &= \sum_{i=1}^{n_{j+1}} (\hat{X}_{j+1})_i = \frac{1}{2} \sum_{i=1}^{n_j/2} (\hat{X}_j)_{2i-1} + (\hat{X}_j)_{2i} \\ &= \frac{1}{2} \left[ (\hat{X}_j)_1 + (\hat{X}_j)_2 + \cdots + (\hat{X}_j)_{n_j} \right] = \underbrace{\frac{n_j}{2}}_{=n_{j+1}} \bar{X}_j = n_{j+1} \bar{X}_j.\end{aligned}\tag{25}$$

By repeated use of this equation we get  $V(\bar{X}_i) = V(\bar{X}_0) = V(\bar{X})$  for all  $0 \leq i \leq d-1$ . This has the consequence that

$$V(\bar{X}) = \frac{\sigma_k^2}{n_k} + e_k \quad \text{for all } 0 \leq k \leq d-1. \tag{26}$$

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 (Master of Science thesis by Marius Jonsson, UiO 2018). It means we can apply blocking transformations until  $e_k$  is sufficiently small, and then estimate  $V(\bar{X})$  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).