

# Statistics and Algebra Reference Guide

Serge Iovleff

March 28, 2019

## Abstract

This reference guide gives a general review of the capabilities offered by the `STK++` library. The library is divided in various *projects*. The "Arrays" project is described in detail in the vignette "STK++ Arrays, User Guide" ([1]) and the quick reference guide. This vignette focus on the "STatistiK" project (which provides statistical tools) and the "Algebra" project (which provide, mainly, matrix decomposition/inversion tools).

## Contents

<b>1</b>	<b>Statistical functors, methods and functions (STatistiK project)</b>	<b>1</b>
1.1	Probabilities . . . . .	1
1.2	Statistical Methods and global functions . . . . .	3
1.2.1	Methods . . . . .	3
1.2.2	Statistical functions . . . . .	4
1.3	Miscellaneous statistical functions . . . . .	5
<b>2</b>	<b>Computing factors</b>	<b>6</b>
<b>3</b>	<b>Linear Algebra classes, methods and functions (Algebra project)</b>	<b>7</b>
3.1	Matrix decomposition . . . . .	7
3.1.1	QR decomposition . . . . .	7
3.1.2	SVD decomposition . . . . .	8
3.1.3	Eigenvalues decomposition . . . . .	9
3.2	Solving least square problems . . . . .	9
3.3	Inverting matrices . . . . .	10

## 1 Statistical functors, methods and functions (STatistiK project)

This section describe the main features provided by the STatistiK project. Mainly

1. the probability classes (1.1),
2. the descriptive statistical methods,
3. the utilities related methods.

The computation of factors using as input a vector or a matrix is detailed in section (2).

### 1.1 Probabilities

All the probabilities handled by R are available in `rtkore`. In the stand-alone `STK++` library, only a subset of theses probabilities are implemented. Probability distribution classes are defined in the `namespace Law` and can be used as in this example

Listing 1: Example

```
#include "STKpp.h"
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    Law::Normal l(2, 1);
    stk_cout << "l.pdf(2)= " << l.pdf(2) << "\n";
    stk_cout << "l.lpdf(2)= " << l.lpdf(2) << "\n";
    stk_cout << "l.cdf(3.96)= " << l.cdf(3.96) << "\n";
    stk_cout << "l.cdfc(3.96)= " << l.cdfc(3.96) << "\n";
    stk_cout << "l.lcdf(3.96)= " << l.lcdf(3.96) << "\n";
    stk_cout << "l.lcdfc(3.96)= " << l.lcdfc(3.96) << "\n";
    stk_cout << "l.icdf(0.975)= " << l.icdf(0.975) << "\n";
    stk_cout << "l.rand()= " << l.rand() << "\n";
    CArray<Real, 2, 3> a; a = 0.5;
    stk_cout << "-----\n";
    stk_cout << "a=\n" << a;
    stk_cout << "a.pdf(1)=\n" << a.pdf(1);
    stk_cout << "a.lpdf(1)=\n" << a.lpdf(1);
    stk_cout << "a.cdf(1)=\n" << a.cdf(1);
    stk_cout << "a.lcdf(1)=\n" << a.lcdf(1);
    stk_cout << "a.cdfc(1)=\n" << a.cdfc(1);
    stk_cout << "a.lcdfc(1)=\n" << a.lcdfc(1);
    stk_cout << "a.icdf(1)=\n" << a.icdf(1);
}
```

Listing 1: Output

```
l.pdf(2) = 0.398942
l.lpdf(2) = -0.918939
l.cdf(3.96) = 0.975002
l.cdfc(3.96) = 0.0249979
l.lcdf(3.96) = -3.68896
l.lcdfc(3.96) = -3.68896
l.icdf(0.975) = 3.95996
l.rand() = 2.89849
-----
a=
0.5 0.5 0.5
0.5 0.5 0.5
a.pdf(1)=
0.129518 0.129518 0.129518
0.129518 0.129518 0.129518
a.lpdf(1)=
-2.04394 -2.04394 -2.04394
-2.04394 -2.04394 -2.04394
a.cdf(1)=
0.0668072 0.0668072 0.0668072
0.0668072 0.0668072 0.0668072
a.lcdf(1)=
-0.0691435 -0.0691435 -0.0691435
-0.0691435 -0.0691435 -0.0691435
a.cdfc(1)=
0.933193 0.933193 0.933193
0.933193 0.933193 0.933193
a.lcdfc(1)=
-0.0691435 -0.0691435 -0.0691435
-0.0691435 -0.0691435 -0.0691435
a.icdf(1)=
2 2 2
2 2 2
```

All probability distribution classes have a similar prototype like the one given below

Listing 2: Prototype of probability distribution class (example taken from Cauchy class)

```
class Cauchy: public IUnivLaw<Real>
{
public:
    Cauchy( Real const& mu=0, Real const& scale=1);
    virtual ~Cauchy() {}
    Real const& mu() const;
    Real const& scale() const;
    void setMu( Real const& mu);
    void setScale( Real const& scale);
    virtual Real rand() const; // generate a Cauchy random variate
    virtual Real pdf( Real const& x) const; // probability distribution function (pdf)
    virtual Real lpdf( Real const& x) const; // log-pdf
    virtual Real cdf( Real const& t) const; // cumulative distribution function (cdf)
    virtual Real cdfc( Real const& t) const; // complementary cdf
    virtual Real lcdf( Real const& t) const; // log-cdf
    virtual Real lcdfc( Real const& t) const; // log-complementary cdf
    virtual Real icdf( Real const& p) const; // inverse cdf (quantiles)
    static Real rand( Real const& mu, Real const& scale);
    static Real pdf( Real const& x, Real const& mu, Real const& scale);
    static Real lpdf( Real const& x, Real const& mu, Real const& scale);
    static Real cdf( Real const& t, Real const& mu, Real const& scale);
    static Real icdf( Real const& p, Real const& mu, Real const& scale);
protected:
    Real mu_;
    Real scale_;
}
```

If  $f$  denote the density of some probability distribution function (pdf) on  $\mathbb{R}$ , the methods have the following meaning

1. `pdf(x)` return the pdf value  $f(x)$ ,
2. `lpdf(x)` return the log-pdf value  $\log(f(x))$ ,
3. `rand()` return a random variate with pdf  $f$ ,
4. `cdf(t)` return the cumulative distribution function (cdf) value  $F(t) = \int_{-\infty}^t f(x)dx$
5. `lcdf(t)` return the log-cdf value  $\log F(t)$
6. `cdfc(t)` return the complementary cdf value  $G(t) = \int_t^{+\infty} f(x)dx$

7. `cdfc(t)` return the log-complementary cdf value  $\log G(t)$
8. `icdf(p)` return the inverse cumulative distribution function value  $F^{-1}(p)$ .

The table 1 gives the list of available probability distribution.

Name	Constructor	R functions	Notes
Bernoulli	<code>Law::Bernoulli(p)</code>	-	
Binomial	<code>Law::Binomial(n,p)</code>	<code>*binom</code>	
Beta	<code>Law::Beta(alpha,beta)</code>	<code>*beta</code>	
Categorical	<code>Law::Categorical(p)</code>	-	p can be any STK++ vector
Cauchy	<code>Law::Cauchy(m,s)</code>	<code>*cauchy</code>	
ChiSquared	<code>Law::ChiSquared(n)</code>	<code>*chisq</code>	
Exponential	<code>Law::Exponential(lambda)</code>	<code>*exp</code>	Parameterization $\lambda e^{-\lambda x}$
FisherSnedecor	<code>Law::FisherSnedecor(df1,df2)</code>	<code>*f</code>	
Gamma	<code>Law::Gamma(a,b)</code>	<code>*gamma</code>	Parameterization $\frac{x^{a-1}}{\beta^a \Gamma(a)} e^{-x/\beta}$
Geometric	<code>Law::Geometric(p)</code>	<code>*geom</code>	
HyperGeometric	<code>Law::HyperGeometric(m,n,k)</code>	<code>*hyper</code>	
Logistic	<code>Law::Logistic(mu,scale)</code>	<code>*logis</code>	
LogNormal	<code>Law::LogNormal(mulog,sdlog)</code>	<code>*lnorm</code>	
NegativeBinomial	<code>Law::NegativeBinomial(size,prob,mu)</code>	<code>*nbinom</code>	
Normal	<code>Law::Normal(mu,sigma)</code>	<code>*norm</code>	
Poisson	<code>Law::Poisson(lambda)</code>	<code>*poiss</code>	
Student	<code>Law::Student(df)</code>	<code>*t</code>	
Uniform	<code>Law::Uniform(a,b)</code>	<code>*unif</code>	
UniformDiscrete	<code>Law::UniformDiscrete(a,b)</code>	-	
Weibull	<code>Law::Weibull(a)</code>	<code>*weibull</code>	

Table 1: List of the available probability distribution in `rtkore`

All Distribution laws methods can be applied to a vector/array/expression using the corresponding methods. An example is given below.

**Listing 3: Compute the log-complementary cdf in a logistic regression**

```
// logis is logistic distribution
STK::Law::Logistic logis;
// y is a (R) matrix of size (n,p) and beta is a (R) vector of size p
STK::RMatrix<double> y;
STK::RVector<double> beta;
// compute log-complementary cdf
(y.*beta).lcdfc(logis);
```

## 1.2 Statistical Methods and global functions

STK++ provides a lot of methods, functions and functors in order to compute usual statistics.

### 1.2.1 Methods

Let  $m$  be any kind of array (square, vector, point, etc...). it is possible to compute the min, max, mean, variance of the elements. These computations can be safe (i.e. discarding N.A. and infinite values) or unsafe and weighted

Method	weighed version	safe versions	Notes
<code>m.norm()</code>	<code>m.wnorm(w)</code>	<code>m.normSafe()</code> ; <code>m.wnormSafe(w)</code>	$\sqrt{\sum  m_{ij} ^2}$
<code>m.norm2()</code>	<code>m.wnorm2(w)</code>	<code>m.norm2Safe()</code> ; <code>m.wnorm2Safe(w)</code>	$\sum  m_{ij} ^2$
<code>m.normInf()</code>	<code>m.wnormInf(w)</code>	<code>m.normInfSafe()</code> ; <code>m.wnormInfSafe(w)</code>	$\sup  m_{ij} $
<code>m.sum()</code>	<code>m.wsum(w)</code>	<code>m.sumSafe()</code> ; <code>m.wsumSafe(w)</code>	$\sum m_{ij}$
<code>m.mean()</code>	<code>m.wmean(w)</code>	<code>m.meanSafe()</code> ; <code>m.wmeanSafe(w)</code>	$\frac{1}{n} \sum m_{ij}$
<code>m.variance()</code>	<code>m.wvariance(w)</code>	<code>m.varianceSafe()</code> ; <code>m.wvarianceSafe(w)</code>	$\frac{1}{n} \sum (m_{ij} - \bar{m})^2$
<code>m.variance(mu)</code>	<code>m.wvariance(mu, w)</code>	<code>m.varianceSafe(mu)</code> ; <code>m.wvarianceSafe(mu, w)</code>	$\frac{1}{n} \sum (m_{ij} - \mu)^2$

Table 2: List of the available statistical methods for the arrays.  $n$  represents the number of elements of  $m$ . For safe versions,  $n$  represents the number of available observations in  $m$ .

### 1.2.2 Statistical functions

For two dimensional arrays, there exists global functions allowing to compute the usual statistics by column or by row. By default all global functions are computing the statistics columns by columns. For example, if  $m$  is an array, `sum(m)` return a row-vector of range `m.cols()` containing the sum of each columns. The alias `sumByCol(m)` can also be used. The sum of each rows can be obtained using the function `sumByow(m)` which return an array of range `m.rows()`.

These computations can be safe (i.e. discarding N.A. and infinite values) or unsafe and/or weighted.

Function	weighed version	safe versions (/ * / for optional arg)
<code>min(m)</code>	<code>min(m, w)</code>	<code>minSafe(m/*, w*/)</code>
<code>minByCol(m)</code>	<code>minByCol(m, w)</code>	<code>minSafeByCol(m/*, w*/)</code>
<code>minByRow(m)</code>	<code>minByRow(m, w)</code>	<code>minSafeByRow(m/*, w*/)</code>
<code>max(m)</code>	<code>max(m, w)</code>	<code>maxSafe(m/*, w*/)</code>
<code>maxByCol(m)</code>	<code>maxByCol(m, w)</code>	<code>maxSafeByCol(m/*, w*/)</code>
<code>maxByRow(m)</code>	<code>maxByRow(m, w)</code>	<code>maxSafeByRow(m/*, w*/)</code>
<code>sum(m)</code>	<code>sum(m, w)</code>	<code>sumSafe(m/*, w*/)</code>
<code>sumByCol(m)</code>	<code>sumByCol(m, w)</code>	<code>sumSafeByCol(m/*, w*/)</code>
<code>sumByRow(m)</code>	<code>sumByRow(m, w)</code>	<code>sumSafeByRow(m/*, w*/)</code>
<code>mean(m)</code>	<code>mean(m, w)</code>	<code>meanSafe(m/*, w*/)</code>
<code>meanByCol(m)</code>	<code>meanByCol(m, w)</code>	<code>meanSafeByCol(m/*, w*/)</code>
<code>meanByRow(m)</code>	<code>meanByRow(m, w)</code>	<code>meanSafeByRow(m/*, w*/)</code>
<code>variance(m, unbiased)</code>	<code>variance(m, w, unbiased)</code>	<code>varianceSafe(m/*, w*/, unbiased)</code>
<code>varianceByCol(m, unbiased)</code>	<code>varianceByCol(m, w, unbiased)</code>	<code>varianceSafeByCol(m/*, w*/, unbiased)</code>
<code>varianceByRow(m, unbiased)</code>	<code>varianceByRow(m, w, unbiased)</code>	<code>varianceSafeByRow(m/*, w*/, unbiased)</code>
<code>varianceWithFixedMean(m, mu, unbiased)</code>	<code>variance*(m, w, mu, unbiased)</code>	<code>variance*Safe(m/*, w*/, mu, unbiased)</code>
<code>varianceWithFixedMeanByCol(m, mu, unbiased)</code>	<code>variance*ByCol(m, w, mu, unbiased)</code>	<code>variance*SafeByCol(m/*, w*/, mu, unbiased)</code>
<code>varianceWithFixedMeanByRow(m, mu, unbiased)</code>	<code>variance*ByRow(m, w, mu, unbiased)</code>	<code>variance*SafeByRow(m/*, w*/, mu, unbiased)</code>

Table 3: List of the available global statistical functions for arrays.  $m$  is the array,  $w$  the vector of weights. `unbiased` is a Boolean to set to `true` if unbiased variance (divided by  $n - 1$ ) is desired.

The covariance can be computed in two ways : using two vectors of same size, or using all the columns (rows) of a two-dimensional array. In the first case the functions return the value of the covariance, in the second case, they return a `CSquareArray`.

Function	weigthed version	safe versions
<code>covariance(v1, v2, unbiased)</code>	<code>covariance(v1, v2, w, unbiased)</code>	<code>covarianceSafe(v1, v2, unbiased)</code> <code>covarianceSafe(v1, v2, w, unbiased)</code>
<code>covarianceWithFixedMean(v1, v2, mean, unbiased)</code>	<code>covariance*(v1, v2, w, mean, unbiased)</code>	<code>covariance*Safe(v1, v2, unbiased)</code>
<code>covariance(m, unbiased)</code> <code>covarianceByRow(m, unbiased)</code>	<code>covariance(m, w, unbiased)</code> <code>covarianceByRow(m, w, unbiased)</code>	
<code>covarianceWithFixedMean(m, mean, unbiased)</code> <code>covarianceWithFixedMeanByRow(m, mean, unbiased)</code>	<code>covariance*(m, w, mean, unbiased)</code> <code>covariance*ByRow(m, w, mean, unbiased)</code>	

Table 4: List of the available covariance functions for vectors and arrays.  $v_1$  and  $v_2$  are vectors,  $m$  is an array,  $w$  a vector of weights. `unbiased` is a Boolean to set to `true` if unbiased covariance (divided by  $n - 1$ ) is desired. The first set of covariance functions return a scalar, the second set of covariance functions return a `CSquareArray`

The following example illustrate the use of the covariance function:

Listing 4: Example

```
#include "STKpp.h"
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    // create coovariance matrix and its Cholesky decomposition
    CArraySquare<Real, 3> s; s << 2.0, 0.8, 0.36,
                                0.8, 2.0, 0.8,
                                0.36,0.8, 1.0;

    Array2DLowerTriangular<Real> L; Array2DDiagonal<Real> D;
    cholesky(s, D, L);
    // create correlated data set
    CArray<Real, 100, 3> a;
    a = a.randGauss() * D.sqrt() * L.transpose() + 1;
    stk_cout << "True sigma=\n" << s;
    stk_cout << "Estimated sigma=\n" << Stat::covariance(a);
    return 0;
}
```

Listing 4: Output

```
True sigma=
 2 0.8 0.36
 0.8 2 0.8
 0.36 0.8 1
Estimated sigma=
 2.48966 1.24486 0.437199
 1.24486 1.95422 0.731717
 0.437199 0.731717 1.0612
```

### 1.3 Miscellaneous statistical functions

Given an array  $m$ , it is possible to center it, to standardize it and to perform the reverse operations. They are listed in table (5) given below. These methods are illustrated in the following example

Listing 5: Example

```
#include "STKpp.h"
using namespace STK;
int main(int argc, char** argv)
{
    CArray<Real, 5, 8> A;
    CArrayPoint<Real, 8> mu, std, mean;
    Law::Normal law(1,2);
    A.rand(law);
    stk_cout << "mean(A)=\n" << (mu=Stat::mean(A));
    stk_cout << "variance(A)=\n"
              << Stat::varianceWithFixedMean(A, mu, false)<< "\n";
    // standardize using empirical mean and standard deviation (std
    // is computed during)
    Stat::standardize(A, mu, std);
    stk_cout << "mean(A)=\n" << (mean=Stat::mean(A));
    stk_cout << "variance(A)=\n"
              << Stat::varianceWithFixedMean(A, mean, false)<< "\n";
    // undo standardization
    Stat::unstandardize(A, mu, std);
    stk_cout << "mean(A)=\n" << (mean=Stat::mean(A));
    stk_cout << "variance(A)=\n"
              << Stat::varianceWithFixedMean(A, mean, false);

    return 0;
}
```

Listing 5: Output

```
mean(A)=
0.183092 1.57987 1.65149 0.816282
 1.97431 1.66366 1.97361
 0.761082
variance(A)=
1.37938 3.69946 3.08708 0.661574
 3.03347 9.26716 3.60422
 12.7621
mean(A)=
-1.11022e-17 -1.33227e-16 2.22045
e-17 6.66134e-17 1.77636e-16
 1.11022e-17 -3.33067e-17
 4.44089e-17
variance(A)=
1 1 1 1 1 1 1 1
mean(A)=
0.183092 1.57987 1.65149 0.816282
 1.97431 1.66366 1.97361
 0.761082
variance(A)=
1.37938 3.69946 3.08708 0.661574
 3.03347 9.26716 3.60422
 12.7621
```

Function	weighted version
<code>center(m, mean)</code>	<code>center(m, w, mean)</code>
<code>centerByCol(m, mean)</code>	<code>centerByCol(m, w, mean)</code>
<code>centerByRow(m, mean)</code>	<code>centerByRow(m, w, mean)</code>
<code>standardize(m, std, unbiased)</code>	<code>standardize(m, w, std, unbiased)</code>
<code>standardizeByCol(m, std, unbiased)</code>	<code>standardizeByCol(m, w, std, unbiased)</code>
<code>standardizeByRow(m, std, unbiased)</code>	<code>standardizeByRow(m, w, std, unbiased)</code>
<code>standardize(m, mean, std, unbiased)</code>	<code>standardize(m, w, mean, std, unbiased)</code>
<code>standardizeByCol(m, mean, std, unbiased)</code>	<code>standardizeByCol(m, w, mean, std, unbiased)</code>
<code>standardizeByRow(m, mean, std, unbiased)</code>	<code>standardizeByRow(m, w, mean, std, unbiased)</code>
<code>uncenter(m, mean)</code>	
<code>uncenterByCol(m, mean)</code>	
<code>uncenterByRow(m, mean)</code>	
<code>unstandardize(m, std)</code>	
<code>unstandardizeByCol(m, std)</code>	
<code>unstandardizeByRow(m, std)</code>	
<code>unstandardize(m, mean, std)</code>	
<code>unstandardizeByCol(m, mean, std)</code>	
<code>unstandardizeByRow(m, mean, std)</code>	

Table 5: List of the available utilities functions for centering and/or standardize an array.  $m$  is an array,  $w$  a vector of weights. If used on the columns, `mean` and `std` have to be points (row-vectors), if used by rows, `mean` and `std` have to be vectors. `unbiased` is a Boolean to set to `true` if unbiased covariance (divided by  $n - 1$ ) is desired.

**Note:**

By default, all functions applied on an array are applied column by column.

## 2 Computing factors

Given a finite collection of object in a vector or an array/expression, it is possible to encode it as factor using the classes `Stat::Factor` (for vectors) and `Stat::MultiFactor` (for arrays). These classes are runners and you have to use the `run` method in order to trigger the computation of the factors.

An example is given below

Listing 6: Example

```
#include "STKpp.h"
using namespace STK;
int main(int argc, char** argv)
{
    CArray<char, 13, 3> fac;
    fac << 'b','a','a', 'c','b','a', 'b','a','a', 'c','a','a', 'd','a','a'
        , 'b','b','a', 'd','c','b', 'c','c','b', 'b','c','b', 'd','b','b'
        , 'b','a','b', 'd','a','b', 'c','c','b';

    Stat::Factor<CArrayVector<char, 13> > f1d(fac.col(1));
    f1d.run();
    stk_cout << "nbLevels= " << f1d.nbLevels() << "\n";
    stk_cout << "asInteger=\n"<< f1d.asInteger().transpose() << "\n";
    stk_cout << "Levels: " << f1d.levels().transpose();
    stk_cout << "Levels counts: " << f1d.counts().transpose();

    Stat::MultiFactor<CArray<char, 13, 3> > f2d(fac);
    f2d.run();
    stk_cout << "nbLevels= " << f2d.nbLevels() << "\n";
    stk_cout << "asInteger=\n"<< f2d.asInteger().transpose() << "\n";
    for (int i=f2d.levels().begin(); i<f2d.levels().end(); ++i)
        stk_cout << "Levels " << i << ": " << f2d.levels()[i].transpose();
    stk_cout << "\n";
    for (int i=f2d.levels().begin(); i<f2d.levels().end(); ++i)
        stk_cout << "Counts " << i << ": " << f2d.counts()[i].transpose();
    return 0;
}
```

Listing 6: Output

```
nbLevels= 3
asInteger=
0 1 0 0 0 1 2 2 2 1 0 0 2

Levels: a b c
Levels counts: 6 3 4
nbLevels= 3 3 2

asInteger=
0 1 0 1 2 0 2 1 0 2 0 2 1
0 1 0 0 0 1 2 2 2 1 0 0 2
0 0 0 0 0 0 1 1 1 1 1 1 1

Levels 0: b c d
Levels 1: a b c
Levels 2: a b

Counts 0: 5 4 4
Counts 1: 6 3 4
Counts 2: 6 7
```

### 3 Linear Algebra classes, methods and functions (Algebra project)

STK++ basic linear operation as product, dot product, sum, multiplication by a scalar,... are encoded in template expressions and optimized.

Since STK++ version 0.9 and later, `lapack` library can be used as back-ends for dense matrix matrix decomposition (QR, Svd, eigenvalues) and least square regression. In order to use `lapack`, you must activate its usage by defining the following macros `-DSTKUSElapack` at compilation time and by linking your code with your installed `lapack` library using `-llapack` (at least in \*nix operating systems).

Class	constructor	Note
lapack::Qr (or Qr)	Qr( data, ref = false) Qr( data)	if data is an <code>ArrayXX</code> and <code>ref</code> is true, data will be overwritten by <code>Q</code>
lapack::Svd (or Svd)	Svd( data, ref = false, withU=true, withV=true) Svd(data)	if <code>ref</code> is true, data will be overwritten by <code>Q</code>
lapack::SymEigen (or SymEigen)	SymEigen( data, ref = false) SymEigen( data)	if data is a <code>SquareArray</code> and <code>ref</code> is true, data will be overwritten by <code>Q</code>
lapack::MultiLeastSquare (or MultiLeastSquare)	MultiLeastSquare( b, a, isBref = false, isAref=false) MultiLeastSquare( b, a)	

#### 3.1 Matrix decomposition

All methods for matrix decomposition are enclosed in classes. Computations are launched using the `run` method.

##### 3.1.1 QR decomposition

QR decomposition (also called a QR factorization) of a matrix is a decomposition of a matrix  $A$  into a product  $A = QR$  of an orthogonal matrix  $Q$  and an upper triangular matrix  $R$ .

QR decomposition can be achieved using either the class `STK::Qr` or if `lapack` is available `STK::lapack::Qr`. In later case, code have to be compiled using `-DSTKUSElapack` flag and linked using `-llapack` (at least for GNU-like compilers).

Listing 7: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    Array2D<Real> a(5,4);
    a << 0, 1, 2, 3,
        2, 3, 4, 5,
        2, 1, 6, 7,
        0, 3,-1, 2,
        3,-1, 1, 1;
    stk_cout << _T("STK++ QR decomposition:\n");
    stk_cout << _T("-----\n");
    Qr q(a); q.run();
    stk_cout << _T("R matrix:\n");
    stk_cout << q.R();
    ArrayXX QR = q.R();
    applyLeftHouseholderArray(QR, q.Q());
    stk_cout << _T("QR matrix:\n");
    stk_cout << QR;
    stk_cout << _T("\nlapack QR decomposition:\n");
    stk_cout << _T("-----\n");
    lapack::Qr lq(a); lq.run();
    stk_cout << _T("R matrix:\n");
    stk_cout << lq.R();
    QR = lq.R();
    applyLeftHouseholderArray(QR, lq.Q());
    stk_cout << _T("QR matrix:\n");
    stk_cout << QR;
    return 0;
}
```

Listing 7: Output

```
STK++ QR decomposition:
-----
R matrix:
4.12311 1.21268 5.57832 6.54846
0 4.41921 2.08981 4.99158
0 0 4.745 4.22321
0 0 0 -1.53827
0 0 0 0

QR matrix:
0 1 2 3
2 3 4 5
2 1 6 7
0 3 -1 2
3 -1 1 1

lapack QR decomposition:
-----
R matrix:
-4.12311 -1.21268 -5.57832 -6.54846
0 -4.41921 -2.08981 -4.99158
0 0 -4.745 -4.22321
0 0 0 -1.53827
0 0 0 0

QR matrix:
0 1 2 3
2 3 4 5
2 1 6 7
0 3 -1 2
3 -1 1 1
```

#### Note:

By default the matrix  $Q$  is represented as a product of elementary reflectors  $Q = H_1 H_2 \dots H_k$ , where  $k = \min(m, n)$  each  $H_i$  has the form  $H_i = I - \tau v v'$ . It is possible to get the  $Q$  matrix (of

size  $(m, m)$  by using the `compQ()` method.  $Q$  will be overwritten.

It is possible to update a QR decomposition when a column is added or removed to the original matrix. In the following example, we remove the column number 2 of a matrix and then insert a column with value 1 after the column number 1.

Listing 8: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    ArrayXX a(5,4), QR;
    a << 0, 1, 2, 3,
        2, 3, 4, 5,
        2, 1, 6, 7,
        0, 3,-1, 2,
        3,-1, 1, 1;
    stk_cout << "lapack QR decomposition:\n";
    lapack::Qr lq(a); lq.run();
    // remove column
    lq.eraseCol(2);
    stk_cout << "R matrix:\n";
    stk_cout << lq.R();
    QR = lq.Q() * lq.R();
    stk_cout << "QR matrix:\n";
    stk_cout << QR;
    // insert constant column with value 1 in column 2
    lq.insertCol(Const::Vector<Real>(5), 2);
    stk_cout << "\nR matrix:\n";
    stk_cout << lq.R();
    QR = lq.Q() * lq.R();
    stk_cout << "QR matrix:\n";
    stk_cout << QR;
    return 0;
}
```

Listing 8: Output

```
lapack QR decomposition:
R matrix:
-4.12311 -1.21268 -6.54846
      0 -4.41921 -4.99158
      0      0  4.49464
      0      0      0
      0      0      0
QR matrix:
0 1 3
2 3 5
2 1 7
0 3 2
3 -1 1
R matrix:
-4.12311 -1.21268 -1.69775 -6.54846
      0 -4.41921 -1.11811 -4.99158
      0      0  0.931381  1.39707
      0      0      0 -4.272
      0      0      0      0
QR matrix:
0 1 1 3
2 3 1 5
2 1 1 7
0 3 1 2
3 -1 1 1
```

### 3.1.2 SVD decomposition

The singular-value decomposition of an  $(m, n)$  real (or complex) matrix  $M$  is a factorization of the form  $U\Sigma V^*$ , where  $U$  is an  $(m, n)$  real (or complex) unitary matrix,  $\Sigma$  is an  $(m, n)$  rectangular diagonal matrix with non-negative real numbers on the diagonal, and  $V$  is an  $(n, n)$  real (or complex) unitary matrix.

The diagonal entries  $\sigma_i$  of  $\Sigma$  are known as the singular values of  $M$ .

Listing 9: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    ArrayXX a(5,4), usvt;
    a << 0, 1, 2, 3,
        2, 3, 4, 5,
        2, 1, 6, 7,
        0, 3,-1, 2,
        3,-1, 1, 1;
    stk_cout << _T("STK++ Svd decomposition:\n");
    stk_cout << _T("-----\n");
    Svd<ArrayXX> s(a); s.run();
    stk_cout << _T("Singular values:\n");
    stk_cout << s.D();
    stk_cout << _T("\nUDV^T matrix:\n");
    stk_cout << s.U()*s.D()*s.V().transpose();
    stk_cout << _T("\nlapack Svd decomposition:\n");
    stk_cout << _T("-----\n");
    lapack::Svd ls(a); ls.run();
    stk_cout << _T("Singular values:\n");
    stk_cout << ls.D().transpose();
    stk_cout << _T("\nUDV^T matrix:\n");
    stk_cout << ls.U()*ls.D().diagonalize()*ls.V().transpose();
    return 0;
}
```

Listing 9: Output

```
STK++ Svd decomposition:
-----
Singular values:
12.6179      0      0      0
      0  4.17637      0      0
      0      0  2.51817      0
      0      0      0  1.00223
UDV^T matrix:
2.5327e-16  1  2  3
      2  3  4  5
      2  1  6  7
5.13478e-16  3 -1  2
      3 -1  1  1
lapack Svd decomposition:
-----
Singular values:
12.6179  4.17637  2.51817  1.00223
UDV^T matrix:
1.07206e-15  1  2  3
      2  3  4  5
      2  1  6  7
-2.35922e-16  3 -1  2
      3 -1  1  1
```



### Note:

Singular values are stored in a vector in `lapack` method and in a diagonal matrix in `STK++` method. It is also possible to compute only  $U$  and/or  $V$  matrix.

### 3.1.3 Eigenvalues decomposition

Let  $A$  be a square  $(n, n)$  matrix with  $n$  linearly independent eigenvectors,  $q_i$  ( $i = 1, \dots, N$ ). Then  $A$  can be factorized as

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^{-1}$$

where  $Q$  is the square  $(n, n)$  matrix whose  $i$ -th column is the eigenvector  $q_i$  of  $A$  and  $\Lambda$  is the diagonal matrix whose diagonal elements are the corresponding eigenvalues, i.e.,  $\Lambda_{ii} = \lambda_i$ .

If  $A$  is a symmetric matrix then  $Q$  is an orthogonal matrix.

STK provide native and `lapack` interface classes allowing to compute the eigenvalue decomposition of a *symmetric* square matrix.

Listing 10: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    CSquareX a(5);
    a << 0, 1, 2, 3, 4,
        2, 3, 4, 5, 6,
        2, 1, 6, 7, 8,
        0, 3, -1, 2, 3,
        3, -1, 1, 1, 0;
    stk_cout << "STK++ Eigen decomposition:\n";
    stk_cout << "-----\n";
    SymEigen<CSquareX> s(a.upperSymmetrize()); s.run();
    stk_cout << "D matrix:\n";
    stk_cout << s.eigenValues();
    CSquareX QDQt = s.eigenVectors() * s.eigenValues().diagonalize()
        * s.eigenVectors().transpose();
    stk_cout << "QDQ^T matrix:\n";
    stk_cout << QDQt;
    stk_cout << "\nlapack Eigen decomposition:\n";
    stk_cout << "-----\n";
    lapack::SymEigen<CSquareX> ls(a); ls.run();
    stk_cout << "D matrix:\n";
    stk_cout << ls.eigenValues();
    QDQt = ls.eigenVectors() * ls.eigenValues().diagonalize() * ls.
        eigenVectors().transpose();
    stk_cout << "QDQ^T matrix:\n";
    stk_cout << QDQt;
    return 0;
}
```

Listing 10: Output

```
STK++ Eigen decomposition:
-----
D matrix:
 20.8996
 0.340126
 -0.39673
 -1.64329
 -8.19967
QDQ^T matrix:
-2.33147e-15  1  2  3          4
              1  3  4  5          6
              2  4  6  7          8
              3  5  7  2          3
              4  6  8  3  1.64313e-14

lapack Eigen decomposition:
-----
D matrix:
-8.19967
-1.64329
-0.39673
0.340126
20.8996
QDQ^T matrix:
-6.88338e-15  1  2  3          4
              1  3  4  5          6
              2  4  6  7          8
              3  5  7  2          3
              4  6  8  3  -7.10543e-15
```

### Note:

STK++ eigenvalues computation need a full symmetric matrix as input while `lapack` version use only upper part of the input data. It is also possible to use the lower part of the matrix in the `lapack` version.

## 3.2 Solving least square problems

In linear regression, the observations are assumed to be the result of random deviations from an underlying relationship between the dependent variables  $y$  and independent variable  $x$ .

Given a data set

$$\{y_{i1}, \dots, y_{id}, x_{i1}, \dots, x_{ip}\}_{i=1}^n$$

of  $n$  statistical units, a linear regression model assumes that the relationship between the dependent variable  $y$  and the regressors  $x$  is linear. This relationship is modeled through a disturbance term that adds "noise" to the linear relationship between the dependent variable and regressors. Thus the model takes the form

$$\mathbf{y}_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n.$$

Often these  $n$  equations are stacked together and written in matrix notation as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

STK++ provide native and lapack interface classes allowing to solve the least square regression problem.

Listing 11: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    CArrayXX y(1000,3), x(1000,5), beta(5,3);
    Law::Normal l(0,1);
    x.rand(l);
    beta << 0, 1, 2,
           2, 3, 4,
           2, 1, 6,
           0, 3,-1,
           3,-1, 1;
    y = x * beta + CArrayXX(1000, 3).rand(l);
    stk_cout << "STK++ MultiLeastSquare:\n";
    stk_cout << "-----\n";
    STK::MultiLeastSquare<CArrayXX, CArrayXX> ols(y,x); ols.run();
    stk_cout << "beta matrix:\n";
    stk_cout << ols.x();
    stk_cout << "\nlapack MultiLeastSquare:\n";
    stk_cout << "-----\n";
    STK::lapack::MultiLeastSquare<CArrayXX, CArrayXX> ls(y,x); ls.run();
    stk_cout << "beta matrix:\n";
    stk_cout << ls.x();
    return 0;
}
```

Listing 11: Output

```
STK++ MultiLeastSquare:
-----
beta matrix:
  0.0247875  1.02139  1.99183
  1.96609   3.01792  4.0079
  2.04481   1.03293  6.03462
-0.00848336  2.99758 -1.05867
  3.02074  -0.996975  1.00954

lapack MultiLeastSquare:
-----
beta matrix:
  0.0247875  1.02139  1.99183
  1.96609   3.01792  4.0079
  2.04481   1.03293  6.03462
-0.00848336  2.99758 -1.05867
  3.02074  -0.996975  1.00954
```

#### Note:

It is also possible to solve weighted least-square regression problems.

### 3.3 Inverting matrices

Matrices can be inverted using either the templated functor `STK::InvertMatrix` or the templated function `STK::invert`. The first example below deals with general square and/or symmetric matrices.

Listing 12: Example

```
#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    CArray<Real, 4, 4> a(4,4);
    a << 0, 1, 2, 3,
        2, 3, 4, 5,
        2, 1, 6, 7,
        0, 3,-1, 2;
    stk_cout << "Inverse general matrix:\n";
    stk_cout << "-----\n";
    stk_cout << a.invert(a);
    stk_cout << "\nInverse upper-symmetric matrix:\n";
    stk_cout << "-----\n";
    stk_cout << a.upperSymmetrize();
    stk_cout << a.upperSymmetrize()*invert(a.upperSymmetrize());
    stk_cout << "\nInverse lower-symmetric matrix:\n";
    stk_cout << "-----\n";
    stk_cout << a.lowerSymmetrize();
    stk_cout << a.lowerSymmetrize()*invert(a.lowerSymmetrize());
    return 0;
}
```

Listing 12: Output

```
Inverse general matrix:
-----
  1 0 0 0
  0 1 0 0
-8.88178e-16 0 1 0
  0 0 0 1

Inverse upper-symmetric matrix:
-----
0 1 2 3
1 3 4 5
2 4 6 7
3 5 7 2
      1 0          0          0
      0 1 -2.22045e-16 -2.22045e-16
6.66134e-16 0          1 -2.22045e-16
-2.10942e-15 0  1.11022e-16          1

Inverse lower-symmetric matrix:
-----
0 2 2 0
2 3 1 3
2 1 6 -1
0 3 -1 2
1 0 0 0
0 1 0 0
0 0 1 0
0 0 0 1
```

The second example deals with lower and upper triangular arrays.

Listing 13: Example

```

#include <STKpp.h>
using namespace STK;
/** @ingroup tutorial */
int main(int argc, char** argv)
{
    Array2DLowerTriangular<Real> a(5,5);
    a << 1,
        1, 2,
        3, 4, 3,
        4, 5, 6, 6,
        7, 8, 2, 3, 2;
    stk_cout << "Inverse lower-triangular matrix:\n";
    stk_cout << "-----\n";
    stk_cout << a*invert(a);
    stk_cout << "\nInverse upper-triangular matrix:\n";
    stk_cout << "-----\n";
    stk_cout << a.transpose()*invert(a.transpose());
    return 0;
}

```

Listing 13: Output

```

Inverse lower-triangular matrix:
-----
1 0 0 0 0
0 1 0 0 0
0 0 1 0 0
0 0 0 1 0
0 0 0 0 1

Inverse upper-triangular matrix:
-----
1 0 0 -1.11022e-16 0
0 1 0 0 0
0 0 1 0 0
0 0 0 1 0
0 0 0 0 1

```

**Note:**

If the matrix is not invertible, the result provided will be a generalized inverse.

**References**

- [1] Serge Iovleff. *STK++ Arrays, User Guide*, 2016. R package version 1.0.2.