

STK++ Utilities reference guide

Serge Iovleff

January 23, 2018

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 and the "Algebra" project.

Contents

1	Statistical functors, methods and functions (STatistiK project)	1
1.1	Probabilities	1
1.2	Statistical Methods and global functions	2
1.2.1	Methods	3
1.2.2	Statistical functions	3
1.3	Miscellaneous statistical functions	4
1.4	Computing factors	5
2	Linear Algebra classes, methods and functions (Algebra project)	6

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.

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.cdf(3.96)= " << l.cdf(3.96) << _T("\n");
    stk_cout << "l.icdf(0.975)= " << l.icdf(0.975) << _T("\n");
    stk_cout << "l.pdf(2)= " << l.pdf(2) << _T("\n");
    stk_cout << "l.lpdf(2)= " << l.lpdf(2) << _T("\n");
    stk_cout << "l.rand()= " << l.rand() << _T("\n");
    CArray33 a;
    a.rand(1);
    stk_cout << "a=\n" << a << _T("\n");
}
```

Listing 1: Output

```
l.cdf(3.96)= 0.975002
l.icdf(0.975)= 3.95996
l.pdf(2)= 0.398942
l.lpdf(2)= -0.918939
l.rand()= 2.62646
a=
0.272322 0.185485 1.41288
1.65554 1.87354 0.168805
2.46967 -0.266933 0.19432
```

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

```

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;
    virtual Real pdf( Real const& x) const;
    virtual Real lpdf( Real const& x) const;
    virtual Real cdf( Real const& t) const;
    virtual Real icdf( Real const& p) const;
    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 value $f(x)$,
2. `lpdf(x)` return the value $\log(f(x))$,
3. `rand()` return a a realization of the density f ,
4. `cdf(t)` return the cumulative distribution function value $F(t) = \int_{-\infty}^t f(x)dx$
5. `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::Bernouilli(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 $\frac{e^{-x/\lambda}}{\lambda}$
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`

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

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 3: Example	Listing 3: Output
<pre>#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; }</pre>	<pre>True sigma= 2 0.8 0.36 0.8 2 0.8 0.36 0.8 1 Estimated sigma= 1.87474 0.809571 0.291616 0.809571 2.22951 0.890907 0.291616 0.890907 1.10475</pre>

1.3 Miscellaneous statistical functions

Given an array m , it is possible to center it, to standardize it and to perform the reverse operations as illustrated in the following example

Listing 4: Example	Listing 4: Output
<pre>#include "STKpp.h" using namespace STK; int main(int argc, char** argv) { CArray<Real, 5, 8> A; CArrayPoint<Real, 8> mu, std, aux; Law::Normal law(1,2); A.rand(law); stk_cout << _T("mean(A) =") << (mu=Stat::mean(A)); stk_cout << _T("variance(A) =") << Stat::varianceWithFixedMean(A, mu, false) << _T("\n\n"); Stat::standardize(A, mu, std); stk_cout << _T("mean(A) =") << (aux=Stat::mean(A)); stk_cout << _T("variance(A) =") << Stat::varianceWithFixedMean(A, aux, false) << _T("\n\n"); Stat::unstandardize(A, mu, std); stk_cout << _T("mean(A) =") << (aux=Stat::mean(A)); stk_cout << _T("variance(A) =") << Stat::varianceWithFixedMean(A, aux, false) << _T("\n\n"); return 0; }</pre>	<pre>mean(A) =1.17981 2.05352 0.339851 1.25085 0.760685 0.878685 0.722511 0.0667084 variance(A) =1.64512 2.2913 3.95029 3.88707 8.99953 1.1543 1.38419 4.13585 mean(A) =5.55112e-17 -2.22045e-17 -2.22045e-17 0 3.26128e-17 2.22045e-17 -17 0 8.88178e-17 variance(A) =1 1 1 1 1 1 1 1 mean(A) =1.17981 2.05352 0.339851 1.25085 0.760685 0.878685 0.722511 0.0667084 variance(A) =1.64512 2.2913 3.95029 3.88707 8.99953 1.1543 1.38419 4.13585</pre>

Function	weighed version
center(m, mean)	center(m, w, mean)
centerByCol(m, mean)	centerByCol(m, w, mean)
centerByRow(m, mean)	centerByRow(m, w, mean)
standardize(m, std, unbiased)	standardize(m, w, std, unbiased)
standardizeByCol(m, std, unbiased)	standardizeByCol(m, w, std, unbiased)
standardizeByRow(m, std, unbiased)	standardizeByRow(m, w, std, unbiased)
standardize(m, mean, std, unbiased)	standardize(m, w, mean, std, unbiased)
standardizeByCol(m, mean, std, unbiased)	standardizeByCol(m, w, mean, std, unbiased)
standardizeByRow(m, mean, std, unbiased)	standardizeByRow(m, w, mean, std, unbiased)
uncenter(m, mean)	
uncenterByCol(m, mean)	
uncenterByRow(m, mean)	
unstandardize(m, std)	
unstandardizeByCol(m, std)	
unstandardizeByRow(m, std)	
unstandardize(m, mean, std)	
unstandardizeByCol(m, mean, std)	
unstandardizeByRow(m, mean, std)	

Table 5: List of the available utilites 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.

1.4 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 array). These classes are runners and you have to use the **run** method in order to trigger the computation of the factors.

An exemple is given below

Listing 5: Example	Listing 5: Output
<pre> #include "STKpp.h" using namespace STK; int main(int argc, char** argv) { CArray<char, 13, 3> datai; datai << '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(datai.col(1)); f1d.run(); stk_cout << _T("nbLevels= ") << f1d.nbLevels() << _T("\n"); stk_cout << _T("asInteger= ") << f1d.asInteger().transpose() << _T("\n"); stk_cout << _T("Levels: ") << f1d.levels().transpose(); stk_cout << _T("Levels counts: ") << f1d.counts().transpose(); Stat::MultiFactor<CArray<char, 13, 3> > f2d(datai); f2d.run(); stk_cout << _T("nbLevels= ") << f2d.nbLevels() << _T("\n"); stk_cout << _T("asInteger=\n") << f2d.asInteger() << _T("\n"); for (int i=f2d.levels().begin(); i<f2d.levels().end(); ++i) { stk_cout << _T("Levels variable ") << i << _T(": ") << f2d.levels()[i]. transpose();} stk_cout << _T("\n"); for (int i=f2d.levels().begin(); i<f2d.levels().end(); ++i) { stk_cout << _T("Levels counts ") << i << _T(": ") << f2d.counts()[i].transpose ();} stk_cout << _T("\n"); return 0; } </pre>	<pre> 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 0 0 1 1 0 0 0 0 1 0 0 2 0 0 0 1 0 2 2 1 1 2 1 0 2 1 2 1 1 0 0 1 2 0 1 1 2 1 Levels variable 0: b c d Levels variable 1: a b c Levels variable 2: a b Levels counts 0: 5 4 4 Levels counts 1: 6 3 4 Levels counts 2: 6 7 </pre>

2 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 backends for dense matrix matrix decompositions (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	Qr(data, ref = false) Qr(data)	if data is an <code>ArrayXX</code> and <code>ref</code> is true, data will be overwritten by Q
lapack::Svd	Svd(data, ref = false, withU=true, withV=true) Svd(data)	if <code>ref</code> is true, data will be overwritten by Q
lapack::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 Q
lapack::MultiLeastSquare	MultiLeastSquare(b, a, isBref = false, isBref=false) MultiLeastSquare(b, a)	

To be continued...

References

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