

Translation into C++ of the R package “Intervalwise Testing for Functional Data”

PACS project (8 credits)

Tutor: **Dr. Alessia Pini**

Student: **Diego Di Mauro**

May 2018



POLITECNICO
MILANO 1863

Introduction

R package for **Intervalwise T^2 Tests with Functional Data**.

Main function **IWT2** plus minor functions to plot the results.

IWT2 function:

- ▶ coded by Alessia Pini (PhD thesis with Simone Vantini),
- ▶ only a draft to check algorithm correctness,
- ▶ not optimized, a bit messy, with memory wastes,
- ▶ very inefficient and almost unusable in applications.

The goal is to implement **IWT2** in C++ in order to make it more efficient, parallel and in order to interface it with the other functions of the R package.

Description

R version

The inputs accepted by **IWT2** are

- ▶ data1: observations of the first population,
- ▶ data2: observations of the second population,
- ▶ mu: mean difference of the two populations under H_0 ,
- ▶ B: number of Monte Carlo iterations,
- ▶ paired: flag for “paired” or “unpaired” test,
- ▶ dx: domain step size,
- ▶ recycle: flag for using “recycle” (for periodic domains),
- ▶ alt: test type or alternative.

Description

R version

The possible test types are

- ▶ “two.sided”: $H_0: \mu_1 - \mu_2 = \mu_0$ vs $H_1: \mu_1 - \mu_2 \neq \mu_0$,
- ▶ “greater”: $H_0: \mu_1 - \mu_2 = \mu_0$ vs $H_1: \mu_1 - \mu_2 > \mu_0$,
- ▶ “less”: $H_0: \mu_1 - \mu_2 = \mu_0$ vs $H_1: \mu_1 - \mu_2 < \mu_0$.

Search for the domain portions that cause H_0 to be rejected:

1. data pre-processing,
2. pointwise test statistics on the original data,
3. pointwise test statistics on the permuted data,
4. intervalwise test statistics,
5. corrections.

Description

R version

The outputs returned by **IWT2** are

- ▶ `test`: test type,
- ▶ `mu`: mean difference of populations under H_0 (same as input),
- ▶ `unadjusted_pval`: pointwise unadjusted p-values,
- ▶ `adjusted_pval`: pointwise adjusted p-values,
- ▶ `pval_matrix`: intervalwise p-values,
- ▶ `data.eval`: `data1` and `data2` juxtaposed,
- ▶ `ord_labels`: population to which the observations belong.

Description

C++ version

The inputs in C++ are

- ▶ data1, data2: observations of the two populations,
- ▶ mu: mean difference under H_0 ,
- ▶ B: number of Monte Carlo iterations,
- ▶ alt: test type or alternative,
- ▶ maxrow: truncation parameter,
- ▶ paired: flag for “paired” or “unpaired” test,
- ▶ recycle: flag for using “recycle” (for periodic domains),
- ▶ THREADS: number of parallel threads to exploit.

Description

C++ version

The outputs in C++ are only

- ▶ `T0`: vector of test statistics of the original data,
- ▶ `pvalue_point`: pointwise unadjusted p-values,
- ▶ `pvalue_inter`: intervalwise p-values,
- ▶ `pvalue_corr`: pointwise adjusted p-values.

Description

C++ version

Libraries

- ▶ `iostream`, `fstream`
- ▶ `vector`, `string`
- ▶ `Eigen/Dense` (version 3.3.3)
- ▶ `ctime`, `ioomanip`

Macros

- ▶ `INFO` (code flow), `SHOW` (partial results), `TIME` (elapsed times)

Typedefs

- ▶ `MatrixType`: `Array<double, Dynamic, Dynamic, RowMajor>`
- ▶ `VectorType`: `Array<double, Dynamic, 1>`
- ▶ `AlterType`: `std::string`

Description

Data reading

For now, we assume the inputs to be available in the text files

- ▶ `Param.txt` with the dimensions $n1$, $n2$, and p ,
- ▶ `Data1.txt` with the $n1 \times p$ elements of `data1` (by row),
- ▶ `Data2.txt` with the $n2 \times p$ elements of `data2` (by row),
- ▶ `Mean0.txt` with the p elements of `mu`.

Later on, we will also see a method to generate them automatically.

Description

Tilde test

We do not really perform the test

$$H_0 : \mu_1 - \mu_2 = \mu_0 \text{ vs } H_1 : \mu_1 - \mu_2 \neq \mu_0,$$

but we actually perform the **tilde test**

$$H_0 : \tilde{\mu}_1 = \mu_2 \text{ vs } H_1 : \tilde{\mu}_1 \neq \mu_2,$$

where $\tilde{\mu}_1 = \mu_1 - \mu_0$.

- ▶ We only care of the difference $\mu_1 - \mu_2$ w.r.t. μ_0 .
- ▶ Many operations will result much simpler.

Description

Inputs check

The inputs data must satisfy some constraints:

- ▶ `data1` and `data2`:
 - ▶ same number of columns,
 - ▶ same number of rows (only in “paired” tests),
- ▶ `mu`: same number of elements as the number of columns of `data1` and `data2`.

In R, the variables `data1`, `data2`, and `mu` are accepted in many formats, while in C++ they are accepted only as numeric variables. Hence, the step size parameter `dx` is not needed any more.

Description

Inputs check

The input parameters must satisfy some constraints:

- ▶ `alt`: one among “two.sided”, “greater”, and “less”,
- ▶ `B`: positive,
- ▶ `maxrow`: between 0 and $p - 1$ (included),

We prefer *int* instead of *unsigned* for many variables so to

- ▶ avoid annoying warnings (if compiling with `-Wall`),
- ▶ rely on safer exit conditions when cycling on decreasing dummy variables.

Description

T0 computation

Let δ be the column-by-column mean differences vector,

```
for (int j = 0; j < p; j++)  
  delta0(j) = data1.col(j).mean() - data2.col(j).mean();
```

then the T^2 test statistic is

- ▶ $\delta \cdot \delta$ in “two.sided” tests,
- ▶ $\delta^+ \cdot \delta^+$ in “greater” tests,
- ▶ $\delta^- \cdot \delta^-$ in “less” tests.

Description

T0 computation

```
VectorType compute_T2 (const VectorType & delta , int p,
                       const AlterType & alt) {

    VectorType T(delta);

    if (!alt.compare("greater")) {
        for (int j = 0; j < p; j++)
            if (delta(j) < 0)
                T(j) = 0;
    }
    else if (!alt.compare("less")) {
        for (int j = 0; j < p; j++)
            if (delta(j) > 0)
                T(j) = 0;
    }

    return T * T;
}
```

Description

Pointwise p-values computation

Pointwise p-values computed via **“two-pop” permutation tests**:

- ▶ we randomize the labels, i.e., we **virtually** exchange some observations between the two populations,
- ▶ we compute the T^2 statistics under the new configuration,
- ▶ we compare the T^2 obtained with the original T_0 ,
- ▶ we estimate the pointwise p-values.

We create the matrix T_{perm} to store the T^2 statistics:

- ▶ B rows, one for each permutation test,
- ▶ p columns, one for each point of the domain.

Description

Paired tests

The two populations must have the same number of observations.

- ▶ We generate a **random binary sequence** of n elements,
- ▶ we **virtually** exchange the observations between the two populations when the sequence has value 1.

```
// Declare indices vector
std::vector<int> indices(n1);

// Generate random binary sequence
for (int i = 0; i < n1; i++)
    indices[i] = rand() % 2;
```


Description

Unpaired tests

- ▶ We generate a **random sample** of the integers $1 : (n1 + n2)$,
- ▶ we **virtually** reorder the observations of both populations accordingly.

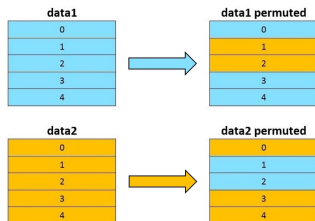
```
// Declare indices vector and auxiliary variables
std::vector<int> indices(n1+n2); bool ok; int k;

// Generate random sample
for (int i = 0; i < (n1+n2); i++) {
    indices[i] = rand() % (n1+n2);
    if (i > 0) {
        ok = false;
        while (ok == false) { // value not chosen yet
            k = 0;
            while (indices[i] != indices[k]) k++;
            if (k < i) indices[i] = rand() % (n1+n2);
            else ok = true;
        }
    }
}
```

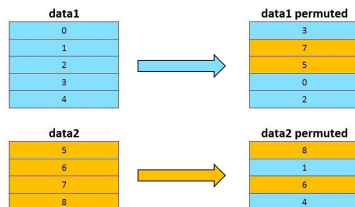
Description

Paired and unpaired tests

`indices = (0, 1, 1, 0, 0)`



`indices = (3, 7, 5, 0, 2, 8, 1, 6, 4)`



Virtually: we do not really exchange the observations in C++. Conversely, in the R version, at each iteration a new copy of data1 and data2 is created and the rows are really exchanged.

Description

Pointwise p-values computation

Pointwise p-values computation:

- ▶ we generate a random sample,
- ▶ we compute the mean differences,
- ▶ we compute the ratio of $T^2 \geq T_0$.

We create the vector count to tally up how many times $T^2 \geq T_0$ and we estimate the p-values with the ratio count/B.

- ▶ Although it is just a counter, we implement count as `std::vector<double>` in order to ease some operations.
- ▶ Besides, even its implementation as `std::vector<unsigned>` brings no performance improvements.

Description

Pointwise p-values computation

```
// Initialize variables
MatrixType T_perm(MatrixType::Zero(B,p));
VectorType count (VectorType::Zero(p));

// Permutation tests
for (int b = 0; b < B; b++) {
    double temp1, temp2;
    VectorType delta(VectorType::Zero(p));

    // Generate random sample...
    // Compute mean differences...

    T_perm.row(b) = compute_T2(delta , p, alt );
    for (int j = 0; j < p; j++)
        if (T_perm(b,j) >= T0(j)) count(j) += 1;
}

VectorType pvalue_point(count/B);
```

Description

Intervalwise p-values computation

Intervalwise p-values computed **for each subinterval**:

- ▶ we sum the T^2 and T_0 of all subinterval points as S^2 and S_0 ,
- ▶ we compare S^2 and S_0 ,
- ▶ we estimate the intervalwise p-values.

We create the matrix `pvalue_inter`:

- ▶ p rows, one for each subinterval cardinality,
(row i : subintervals with $p - i$ points)
- ▶ p columns, one for each subinterval starting point.
(column j : subintervals starting at point j)

Description

Recycle

- ▶ Using “recycle” means assuming a **periodic/cyclic domain**.
- ▶ We have always p subintervals for any length.
- ▶ The matrix `pvalue_inter` is full.

0	0 1 2 3 4 5	1 2 3 4 5 0	2 3 4 5 0 1	3 4 5 0 1 2	4 5 0 1 2 3	5 0 1 2 3 4
1	0 1 2 3 4	1 2 3 4 5	2 3 4 5 0	3 4 5 0 1	4 5 0 1 2	5 0 1 2 3
2	0 1 2 3	1 2 3 4	2 3 4 5	3 4 5 0	4 5 0 1	5 0 1 2
3	0 1 2	1 2 3	2 3 4	3 4 5	4 5 0	5 0 1
4	0 1	1 2	2 3	3 4	4 5	5 0
5	0	1	2	3	4	5
	0	1	2	3	4	5

Description

No recycle

- ▶ Not using “recycle” means assuming an **acyclic domain**.
- ▶ We have less and less subintervals as their length increases.
- ▶ The matrix `pvalue_inter` is lower-triangular.
- ▶ The upper-triangular part is equal to 0 (C++) or NaN (R).

0	0 1 2 3 4 5	_____	_____	_____	_____	_____
1	0 1 2 3 4	1 2 3 4 5	_____	_____	_____	_____
2	0 1 2 3	1 2 3 4	2 3 4 5	_____	_____	_____
3	0 1 2	1 2 3	2 3 4	3 4 5	_____	_____
4	0 1	1 2	2 3	3 4	4 5	_____
5	0	1	2	3	4	5
	0	1	2	3	4	5

Description

Intervalwise p-values computation

The **intervalwise p-values computation** cycles on the domain subintervals from the shortest to the longest ones, hence the matrix `pvalue_inter` is filled upwards.

Nevertheless, if the **truncation** parameter `maxrow` is positive, the cycles stop at the subintervals of cardinality $p - \text{maxrow}$.

Duplication: the R version duplicates the vector `T0` and the matrix `T_perm` to ease some operations, thus producing a significant waste of memory (if $B = 100.000$ and $p = 1.000$, the matrix `T_perm` is almost 1 GB). In C++, we rely only on already existing variables.

Description

Intervalwise p-values computation

```
// Cycle on subinterval length (p - i)
for (int i = p - 2; i >= maxrow; i--) {

    // Cycle on subinterval starting point (j)
    for (int j = 0; j < p; j++) {
        cont = 0;
        // Sum T0...

        // Cycle on permutations
        for (int b = 0; b < B; b++) {
            // Sum T2...

            if (T_temp >= T0_temp) cont++;
        }

        // Compute p-values
        pvalue_inter(i,j) = cont / double(B);
    }
}
```

Description

Corrections

In order to obtain theoretically sound results, we apply some corrections to the p-values computed so far, by combining pointwise and intervalwise p-values.

H_0 should be rejected at a point not only if the pointwise p-value at that point is low, but if also the intervalwise p-values of all subintervals to which that point belongs are low.

Then, the **corrected p-value** at point j is the maximum among the pointwise p-value at point j and the intervalwise p-value of all subintervals containing point.

Description

Corrections with and without recycle

0	0 1 2 3 4 5	1 2 3 4 5 0	2 3 4 5 0 1	3 4 5 0 1 2	4 5 0 1 2 3	5 0 1 2 3 4
1	0 1 2 3 4	1 2 3 4 5	2 3 4 5 0	3 4 5 0 1	4 5 0 1 2	5 0 1 2 3
2	0 1 2 3	1 2 3 4	2 3 4 5	3 4 5 0	4 5 0 1	5 0 1 2
3	0 1 2	1 2 3	2 3 4	3 4 5	4 5 0	5 0 1
4	0 1	1 2	2 3	3 4	4 5	5 0
5	0	1	2	3	4	5
	0	1	2	3	4	5

0	0 1 2 3 4 5	—	—	—	—	—
1	0 1 2 3 4	1 2 3 4 5	—	—	—	—
2	0 1 2 3	1 2 3 4	2 3 4 5	—	—	—
3	0 1 2	1 2 3	2 3 4	3 4 5	—	—
4	0 1	1 2	2 3	3 4	4 5	—
5	0	1	2	3	4	5
	0	1	2	3	4	5

Description

Corrections

Duplication: the R version duplicates the $p \times p$ matrix `pvalue_inter` so to simplify some operations with periodic domains but - in so doing - it causes a significant memory waste. In the C++ version, we avoid that by playing with the indices.

Overwrite: the R version generates a $p \times p$ matrix to perform the corrections but - in the end - it needs only the `maxrow`-th row of it. In C++, we only create a p -dimensional vector and we overwrite it at each iteration, instead of wasting p new elements every time.

Description

Corrections

```
// Cycle on subintervals length
for (int i = maxrow; i < p; i++) {

    if (recycle==true) {

        // Cycle on points
        for (int j = p - 1; j >= 0; j--) {

            if (j==p-1) temp_last = pvalue_corr(p-1);

            if (i==maxrow) pvalue_corr(j) = pvalue_inter(i,j);

            else if (j==0) pvalue_corr(j) = max(pvalue_inter(i,j),
                                                max(pvalue_corr(j),
                                                    temp_last));

            else                pvalue_corr(j) = max(pvalue_inter(i,j),
                                                max(pvalue_corr(j-1),
                                                    pvalue_corr(j)));

        }
    }
    else { } // similar...
}
```

Parallelization

OpenMP is more appropriate for Monte Carlo methods (we repeat B times the same operations) because it works in shared memory. **MPI** would need to broadcast both matrices to all slave nodes.

In the **pointwise computation**, we parallelize the B permutations

- ▶ as equally as possible,
- ▶ “statically” (same complexity for each permutation).

In the **intervalwise computation**, we parallelize the subintervals length

- ▶ “dynamically” (increasing complexity for increasing length).

R interface

```
#include <RcppEigen.h>
#include <omp.h>
// [[Rcpp::depends(RcppEigen)]]
// [[Rcpp::plugins(openmp)]]

// Includes, Macros, Typedefs...

//[[Rcpp::export]]
List IWT2OMP (int B = 1000, AlterType alt = "two.sided",
              int maxrow = 0, bool paired = false,
              bool recycle = false, int THREADS = 1) {
  // Algorithm...

  List RES;
  RES["T0"] = T0;
  RES["point"] = pvalue_point;
  RES["inter"] = pvalue_inter;
  RES["corre"] = pvalue_corr;
  return RES;
}
```

Results

Parameters: n , p

n	p	Time (seconds)
50	100	5
100	100	5
200	100	5
50	200	47
100	200	47
200	200	47
50	500	408
100	500	409
200	500	410

Results

Parameters: maxrow, B

maxrow	Time (seconds)
0	15
20	12
100	5

B	p	Time (seconds)
1000	200	2
5000	200	17
10000	200	62
1000	500	20
5000	500	101
10000	500	406

Results

Parameters: `paired`, `recycle`, `alt`

<code>paired</code>	<code>recycle</code>	Time (seconds)
<code>true</code>	<code>true</code>	20
<code>false</code>	<code>true</code>	20
<code>true</code>	<code>false</code>	7
<code>false</code>	<code>false</code>	7

<code>alt</code>	p	Time (seconds)
<code>"two.sided"</code>	200	15
<code>"greater"</code>	200	15
<code>"less"</code>	200	15
<code>"two.sided"</code>	500	412
<code>"greater"</code>	500	403
<code>"less"</code>	500	407

Results

Sections

The time needed by each section of the algorithm depends on the data and problem dimensions such as $n1$, $n2$, p , B , and maxrow .

- ▶ Read: $\sim 5\%$ (increases with n)
- ▶ T0: $\sim 0\%$ (increases with p)
- ▶ Point: $\sim 1\%$ (increases with p)
- ▶ Interval: $\sim 95\%$ (increases with B , decreases with maxrow)
- ▶ Correct: $\sim 0\%$ (increases with p)

Results

R vs C++

C++ version averagely takes 20 times less than R version...

p	B	recycle	C++ (min)	R (min)	Ratio
500	1000	true	0.33	9	4%
500	5000	true	1.68	48	3%
500	10000	true	6.77	100	7%
200	10000	true	0.28	10	3%
200	10000	false	0.16	2	8%

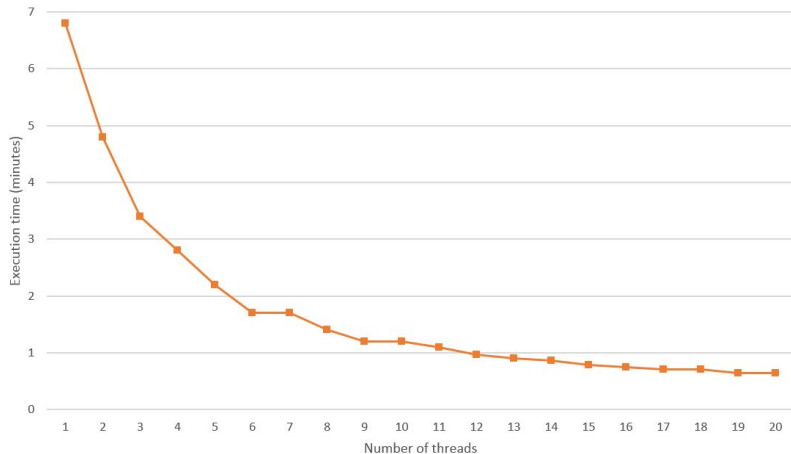
...even without parallelization!

Results

Parallelization

R version: > 3 hours

C++ version: < 2 minutes



Possible extensions

- ▶ **Domain**
 - ▶ 2-dimensional domains
 - ▶ non-equidistant points
- ▶ **Intervalwise p-values computation**
 - ▶ stop when a p-value exceeds a threshold
- ▶ **Test types**
 - ▶ one-population tests
 - ▶ multi-population tests
- ▶ Possibility to set **seed** for random generations

Acknowledgements

I believe that this project will be of great use since it has a lot of applications. Due to its original very long execution times, it was not used much in the past, but now it can be run in less than 1% of the time than before!

I would like to thank everyone who helped me with this project along this long year, in particular,

- ▶ Alessia Pini,
- ▶ Aymeric Stamm,
- ▶ Carlo De Falco,
- ▶ Luca Paglieri.