

Progetto Ingegneria Informatica

Datastream Change Detection

Boracchi Giacomo Antonio Miele Carrera Diego

Fioravanti Massimo

10 luglio 2018

Indice

1	introduzione	2
1.1	Descrizione del problema	2
2	Fondamenti teorici	4
2.1	Statistica	4
2.2	Test di ipotesi	4
3	Implementazione	5
3.1	Implementazione C	5
3.2	Test Multivariati	7
4	Ottimizzazione	9
4.1	Permutazioni a spazio costante e visita ordinata	9
4.2	Campioni di dimensione qualsiasi	10
5	Conclusioni	12
5.1	Correttezza	12
5.2	Lavori Futuri	13

1 introduzione

1.1 Descrizione del problema

Test di ipotesi Varie branche scientifiche e economiche moderne, come ad esempio la diagnostica medica, si basano sulla capacità di discernere le differenze tra popolazioni di oggetti. Se si riesce a descrivere ciò che si sta analizzando tramite variabili numeriche allora ne consegue che si sta affrontando un problema di natura matematica. Le procedure che permettono di operare tali distinzioni sono dette test di verifica di ipotesi e la loro formulazione rigorosa giace nel dominio della statistica. Talvolta il costo computazionale necessario per risolvere tali operazioni è assai elevato, in particolare perché la complessità degli algoritmi è spesso legata alla dimensione del campione. Inoltre se si tenta di analizzare un evento periodico reale che avviene ad alta frequenza è importante riuscire a garantire la maggiore velocità di analisi possibile, poiché ciò permette di estendere l'uso dei test di ipotesi a nuovi ambiti. Per queste ragioni è interessante studiare implementazioni efficienti dei test statistici.

Permutation Test Tal volta è possibile porre delle assunzioni riguardo quali sono i parametri delle popolazioni che si vuole studiare, in tali casi è spesso possibile creare dei test specifici particolarmente potenti. Tali test prendono il nome di test parametrici. Altre volte le caratteristiche delle popolazioni sono completamente ignote, poiché magari derivano da eventi naturali. I test ideati per lavorare con queste popolazioni prendono il nome di test non parametrici. I test di permutazione sono particolari all'interno di questa categoria perché il loro funzionamento si basa sul calcolare una statistica riguardanti le popolazioni che si intende confrontare, per poi permutare casualmente i due gruppi e ripetere il calcolo. Si itera questa procedura sino ad aver raccolto dati a sufficienza, e si confronta la statistica della popolazione originale contro una soglia calcolata dalle statistiche delle permutazioni permutate. Se il primo numero risulta essere maggiore del secondo allora ne consegue che le due popolazioni non possono disporre delle stesse caratteristiche, poiché se quello fosse stato il caso allora non si discosterebbe così tanto dal caso medio ottenuto dalle permutazioni.

Schede Grafiche Gli ultimi 15 anni hanno visto l'affermarsi della tecnologia delle schede grafiche come strumento utile alla computazione ad alta velocità. Progettate inizialmente per la creazione di immagini 3D da presentare a schermo, tali unità di calcolo presentano caratteristiche che le rendono adatte a eseguire calcoli in parallelo. Una scheda grafica è composta svariate unità di calcolo ognuna disposta dei propri registri privati, queste unità sono poi raggruppate a gruppi i quali condividono una memoria locale e infine tutte loro accedono alla stessa memoria globale. I gruppi condividono anche parte delle componenti hardware, in particolare quelle legate al controllo del flusso del programma. Ciò implica che mentre il consumo energetico è minore rispetto alle cpu convenzionali, le unità di calcolo all'interno di un gruppo devono necessariamente eseguire le stesse operazioni. Tale tecnica prende il nome di SIMD, cioè Single Instruction Multiple Data, poiché disponendo di registri diversi per ogni nucleo di calcolo è possibile raggiungere risultati distinti su ognuno di essi, anche se la procedura è la stessa.

Tale architettura permette quindi di accelerare tutti gli algoritmi che possono essere eseguiti in parallelo. Se si seleziona una statistica adatta, questo è esattamente il caso riguardante i permutation test, poiché le permutazioni non necessitano di conoscere quale fosse stato il risultato della permutazione precedente, mentre il calcolo della statistica segue le stesse operazioni in ogni unità di calcolo.

Risultati In questo documento si presenterà una soluzione software del problema sopra descritto destinata all'uso su schede grafiche dotata delle seguenti caratteristiche:

- **Test multivariati** l'analisi di campioni multidimensionali risulta particolarmente difficile e costosa, la nostra implementazione risolve tale problema per costruzione.
- **Memoria lineare** il nostro algoritmo richiede memoria pari a quella necessaria per allocare l'input in ingresso. Tale soluzione è ideale poiché la memoria disponibile sulle schede grafiche è limitata e garantire tale proprietà nel problema affrontato è difficile.
- **Accesso sequenziale** il nostro algoritmo accede alla memoria in maniera sequenziale ed esattamente una volta per ogni elemento, tale soluzione è ideale poiché l'accesso non ordinato rallenta significativamente l'efficienza delle schede grafiche.

Draft

2 Fondamenti teorici

2.1 Statistica

Definiamo come statistica una funzione che associa a una popolazione un numero reale che la rappresenta. Ad esempio la funzione che associa ad una popolazione di lunghezze la media di tali valori è una statistica.

2.2 Test di ipotesi

Per test di verifica di ipotesi si intende un test ideato per verificare su un'ipotesi risulta essere vera o meno, e per ipotesi si intende un'affermazione riguardante oggetti del mondo reale. I test di ipotesi si dividono in test di ipotesi deterministici o statistici, ed è su questi ultimi che ci concentreremo.

L'ipotesi da verificare prende il nome di H_0 o ipotesi nulla. Si definisce invece l'ipotesi contraria come H_1 o ipotesi alternativa. Come conseguenza di un test vi possono essere due tipi di errori, rifiutare l'ipotesi nulla quando essa era vera, come ad esempio non identificare una malattia in un paziente malato, oppure accettare l'ipotesi nulla quando essa è falsa, ad esempio identificare una malattia in un paziente sano. Il primo caso prende il nome di errore di primo tipo, e il secondo errore di secondo tipo. Minimizzare gli errori del primo tipo senza curarsi degli errori del secondo è triviale, poiché è sufficiente accettare ogni volta l'ipotesi. Similmente è vero l'inverso, rifiutare sempre l'ipotesi H_0 permette di minimizzare gli errori del secondo tipo ma garantisce di massimizzare quelli del primo. Di conseguenza la maniera corretta di formulare un test è garantire costanti gli errori del primo tipo, ad esempio il 5% dei casi, e si desidera minimizzare gli errori del secondo tipo. La soglia mantenuta costante è detta α . Il valore calcolato come 1 meno errori del secondo tipo è detto potenza del test.

Come conseguenza del fatto che gli errori del primo tipo devono essere costanti ne consegue che per certe statistiche è possibile calcolare dei valori di soglia che quando confrontati con la statistica valutata per una popolazione che rispetta l'ipotesi h_0 soddisfano la richiesta che gli errori del primo tipo siano esattamente $\alpha\%$. Questi valori, detti soglie, sono poi utilizzati come di valori di confronto nel caso di popolazioni che si vuole determinare se soddisfano l'ipotesi nulla o meno, e in questa maniera si può ricavare una procedura algoritmica che permetta di eseguire il test di ipotesi.

Permutation Test In alcuni casi le soglie dipendono dalla struttura dei dati in ingresso, in particolare quando i dati hanno una natura vettoriale, e di conseguenza non è possibile calcolarle in anticipo. In tali situazioni è necessario affidarsi ad altre soluzioni e una di queste è il permutation test che permette di stabilire se due gruppi provengono dalla stessa popolazione. La considerazione sulla quale si basa un permutation test è che se si permutano i due gruppi casualmente allora in media $\alpha\%$ delle volte la statistica calcolata avrà un valore maggiore della soglia ignota, poiché permutando il campione allora sarà impossibile distinguere gli elementi appartenenti al primo gruppo da quelli appartenenti al secondo. Se si ordinano le statistiche valutate sulle popolazioni permutate e si seleziona quella nella posizione numero permutazioni $\cdot (1-\alpha)$. Tale valore risulta essere una stima della soglia, ed è quindi possibile eseguire test di ipotesi statistica pur essendo del tutto ignari della distribuzione di una popolazione.

3 Implementazione

3.1 Implementazione C

L'algoritmo è composto da due fasi, il calcolo delle statistiche per ogni permutazione e il confronto dei risultati ottenuti. Poiché il risultato deve essere un booleano che indica se le popolazioni dispongono delle stesse caratteristiche, e la statistica è indipendente dal test allora ne consegue la intestazione della stringa `c` può essere la seguente

```
1 int runPermutationTest(Data *data, float (*statistic)(Data* data));
```

`Data` è la struct contenente è definita come segue

```
1
2 typedef struct data
3 {
4     float* sample;
5     unsigned size;
6     unsigned cutPoint;
7     float alpha;
8     float iterations;
9 } Data;
```

dove `sample` è il puntatore al primo elemento dell'array dei dati che include entrambe le popolazioni, `cutPoint` è la dimensione della prima popolazione (e di conseguenza l'indice del primo elemento della seconda) e `size` è la somma della dimensione del primo campione e del secondo. La dimensione del secondo campione può essere ricavata come `size - cutPoint`. `Alpha` è il livello di significatività del test, mentre `iterations` è il numero di iterazioni da effettuare.

Si noti anche il secondo parametro

```
1 float (*statistic)(Data* data)
```

tale scrittura indica che il secondo parametro è una funzione che accetta un puntatore a `Data` e ritorna un `float`, cioè una funzione che rappresenta una statistica.

la funzione diventa trivialmente.

```
1
2 int runPermutationTest(Data *data, float (*statistic)(Data* data))
3 {
4     float statistics[data->iterationsCount];
5
6     //calcolo della statistica delle popolazioni non permutate
7     float sampleStatistic = statistic(data);
8
9     //calcolo delle statistiche permutate
10    for (unsigned a = 0; a < data->iterations; a++)
11    {
12        permute(data);
13        statistics[a] = statistic(data);
14    }
15
16    //confronto della statistica con quelle permutate
17    sort(statistics, data->iterations, sizeof(float));
18    int k = data->iterations - ((int)(data->alpha * data->iterations));
19    return statistics[k] > sampleStatistic;
20
21 }
```

La parte in cui si gestisce il caso \geq del return è stato ignorato in favore della semplicità della descrizione.

OpenCL

Riguardo opencl Si è mostrato come l'implementazione in un linguaggio di programmazione regola risulta essere triviale, discuteremo ora di come implementare tale algoritmo in OpenCL, cioè uno dei linguaggi utilizzati dalle schede grafiche. Un programma OpenCL prende il nome di kernel. Un kernel è fondamentalmente una funzione che viene invocata da un normale programma risiedente nella macchina ospite della scheda. L'invocatore deve specificare quali sono i parametri di invocazione, e quale è l'area memoria dove il kernel può scrivere il risultato computato. Risulta inoltre necessario specificare quante unità di calcolo vogliamo dedicare all'operazione, anche se tale valore non è accessibile all'interno di un programma OpenCL.

Ad esempio supponiamo di voler scrivere un kernel che raddoppia il valore di ogni elemento di un array

```
1 __kernel void double(__global const float *in, __global float *out)
2 {
3     unsigned int i = get_global_id(0);
4
5     out[i] = in[i] * 2;
6 }
```

Il programma invocante caricherà nel vettore in i dati da raddoppiare, creerà l'array out dal quale desidera ricavare i risultati e infine lancerà il kernel, indicando che desidera dedicare un numero di core pari alla dimensione dell'array.

Ogni core utilizzerà la funzione `get_global_id` per ricavare il proprio offset rispetto alla prima unità di calcolo e utilizzerà questo valore per scoprire quale numero deve raddoppiare e dove deve essere salvato.

La keyword `__global` serve per indicare alla scheda che è un riferimento alla memoria condivisa tra tutte le unità di calcolo.

Implementazione permutation test Risulta triviale identificare quale parti dell'implementazione descritta nella sezione 3.1.

I dati originali vengono ospitati sulla scheda, il ciclo contenuto del ciclo for viene spostato nel kernel che si desidera scrivere mentre l'ordinamento e il confronto vengono operati sulla cpu.

Risulta quindi che un kernel è composto dalla permutazione e dalla valutazione della statistica.

```
1 __kernel void p_test
2 (
3     __global const float* in,
4     __global float* out,
5     unsigned int cutPoint,
6     unsigned int sampleSize,
7 )
8 {
9     unsigned int i = get_global_id(0);
10    float[sampleSize] permutated;
11 }
```

```

12     permutate(in, permuted, i);
13
14     b[i] = evaluateStatistic(in, sampleSize, cutPoint, vectorSize);
15 }

```

La scrittura efficiente della statistica verrà analizzata nella sezione ??, quindi non verrà descritta ora. L'intero problema si riduce alla capacità di permutare i dati in ingresso.

Ciò è possibile basandosi sulla seguente osservazione: Sia P un numero primo qualsiasi, allora per tutti i numeri $X, Y < P / 2$ vale che $(X * X) \bmod P \neq (Y * Y) \bmod P$. Ciò significa che iterando attraverso tutti i numeri minori di P è possibile ottenere una sequenza di numeri che non si ripetono mai.

```

1 unsigned int singlePermutate(unsigned int x, unsigned int prime)
2 {
3     if (x <= prime / 2)
4         return x * x % prime;
5     else
6         return prime - (x * x % prime);
7
8 }

```

Si noti che alcune ottimizzazioni sono necessarie per rompere la simmetria centrale dovuta al funzionamento della formula, ma verranno tralasciate per semplicità.

Da questa funzione segue che la variabile X rappresenta il seme del generatore di numeri casuali, e quindi è sufficiente a rappresentarne completamente lo stato.

Il kernel diviene quindi, a patto che sample size sia uguale un numero primo

```

1 __kernel void p_test
2 (
3     __global const float* in,
4     __global float* out,
5     unsigned int cutPoint,
6     unsigned int sampleSize,
7 )
8 {
9     unsigned int i = get_global_id(0);
10    float[sampleSize] permuted;
11
12    for (int a = 0; a < sampleSize; a++)
13        permuted[i] = int[permuteSingle(a, sampleSize)];
14
15    b[i] = evaluateStatistic(in, sampleSize, cutPoint, vectorSize);
16 }

```

Risulta quindi che tutto ciò che rimane da fare è ordinare i risultati e confrontarli con la statistica dei dati non premutati. Ciò può svolto dal programma risiedente nella memoria principale della macchina, poiché è necessario attendere che tutte le statistiche vengano valutate e nel caso si utilizzasse più di una scheda non vi è alcuna garanzia che esse terminino in contemporanea.

3.2 Test Multivariati

Queste ottimizzazioni riguardano esclusivamente la maniera in cui si accede ai dati depositati nella memoria della scheda. Di conseguenza è possibile caricare dati sotto forma di struct, invece che semplici float. Risulta quindi raggiunto

uno degli obbiettivi di tale documento, disponiamo di un implementazione di un test di ipotesi capace di operare su campioni multivariati che non necessita di precalcolare le soglie.

Draft

4 Ottimizzazione

4.1 Permutazioni a spazio costante e visita ordinata

Poiché le computazioni avvengono in parallelo, non è possibile permutare sul posto i dati in ingresso, poiché ciò richiederebbe una quantità di memoria pari alla dimensione dei dati in ingresso moltiplicati per il numero di permutazioni, poiché ogni unità di calcolo deve allocare localmente i dati permutati. Risulta quindi imperativo utilizzare un algoritmo a spazio costante, poiché la memoria disponibile per ogni singolo core è limitata. Si mostrerà ora come è possibile computare statistiche che operano in maniera sequenziale senza precalcolare la permutazione da applicare ai dati.

Prendiamo come esempio la statistica che date due popolazioni ritorna la differenza delle loro medie.

```
1 float evaluateStatistic
2 (
3     __global const float* in ,
4     const unsigned int  sampleSize ,
5     const unsigned int  cutPoint
6 )
7 {
8     float meanSample1 = 0;
9     float meanSample2 = 0;
10
11    for (int a = 0; a < cutPoint; a++)
12        meanSample1 += in[a];
13    meanSample1 /= cutPoint;
14
15    for (int a = cutPoint < sampleSize; a++)
16        meanSample2 += in[a];
17    meanSample2 /= sampleSize - cutPoint;
18
19    return abs(meanSample1 - meanSample2);
20 }
```

grazie alla funzione di permutazione possiamo mappare gli indici sequenziali al loro equivalente permutato, così facendo perdiamo però la possibilità di capire se l'elemento appartiene al primo gruppo o al secondo. La funzione deve quindi essere modificata in:

```
1 float evaluateStatistic
2 (
3     __global const float* in ,
4     int prime ,
5     int index ,
6     const unsigned int  sampleSize ,
7     const unsigned int  cutPoint
8 )
9 {
10    float meanSample1 = 0;
11    float meanSample2 = 0;
12
13    for (int a = 0; a < sampleSize; a++)
14    {
15        int realIndex = permutateSingle(index++, prime);
16
17        if (realIndex < cutPoint)
18            meanSample1 += in[realIndex];
19        else
```

```

20         meanSample2 += in[realIndex];
21     }
22
23     meanSample1 /= cutPoint;
24     meanSample2 /= sampleSize - cutPoint;
25
26     return abs(meanSample1 - meanSample2);
27 }

```

Si noti che questo algoritmo dispone di un'altra caratteristica, l'accesso alla memoria è sequenziale e tutti i nuclei accedono contemporaneamente allo stesso elemento. Poiché le memorie globali delle schede grafiche sono molto più lente che le memorie private, questa soluzione migliora sensibilmente le performance dell'algoritmo.

4.2 Campioni di dimensione qualsiasi

Come precedentemente indicato, per colpa del generatore di numeri casuali siamo costretti a utilizzare campioni la cui dimensione è un numero primo. Tale restrizione è particolarmente insoddisfacente, quindi si discuterà ora di come rimuoverla: Se la dimensione del campione è minore del numero primo è possibile allocare come variabile locale un array la cui dimensione è pari alla differenza tra i due numeri. Ogni volta che si estrae dalla permutazione un numero maggiore della dimensione del campione si salva il valore proveniente dall'array di input nell'array locale in maniera sequenziale. Dopo aver processato tutti gli elementi dell'array globale si processano nuovamente tutti quelli allocati localmente, con la garanzia che ora l'indice estratto sarà minore della dimensione del campione, poiché tutti gli indici maggiori sono già stati valutati.

```

1 float evaluateStatistic
2 (
3     __global const float* in,
4     int prime,
5     int index,
6     const unsigned int sampleSize,
7     const unsigned int cutPoint
8 )
9 {
10     float meanSample1 = 0;
11     float meanSample2 = 0;
12     float[prime - sampleSize] local;
13     int lastAllocatedIndex = 0;
14
15
16     for (int a = 0; a < sampleSize; a++)
17     {
18
19         int realIndex = permutateSingle(index++, prime);
20         if (realIndex > sampleSize)
21             local[lastAllocatedIndex++] = in[realIndex];
22         else if (realIndex < cutPoint)
23             meanSample1 += in[realIndex];
24         else
25             meanSample2 += in[realIndex];
26     }
27
28     for (int a = 0; a < prime - sampleSize; a++)
29     {
30         int realIndex = permutateSingle(index++, prime);

```

```

31
32     if (realIndex > sampleSize)
33         local[lastAllocatedIndex++] = local[realIndex];
34     else if (realIndex < cutPoint)
35         meanSample1 += local[realIndex];
36     else
37         meanSample2 += local[realIndex];
38 }
39
40 meanSample1 /= cutPoint;
41 meanSample2 /= sampleSize - cutPoint;
42
43 return abs(meanSample1 - meanSample2);
44 }

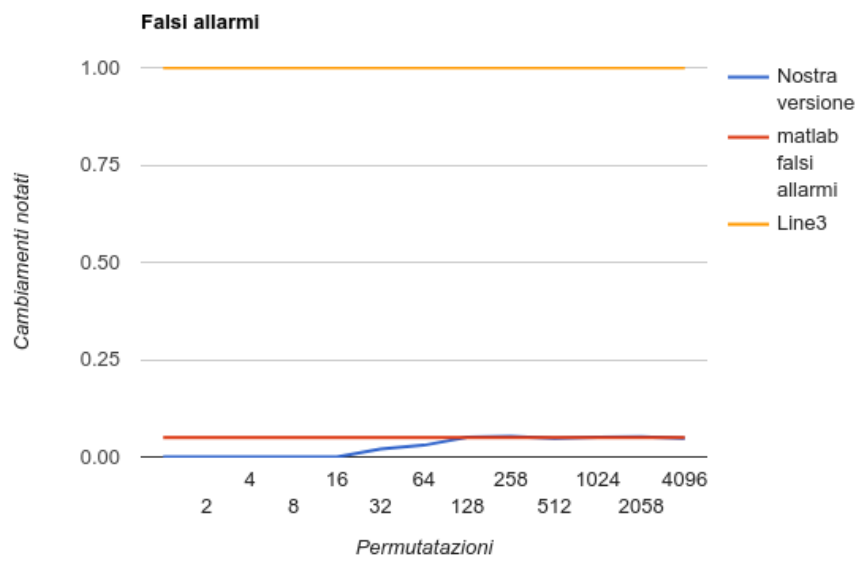
```

Draft

5 Conclusioni

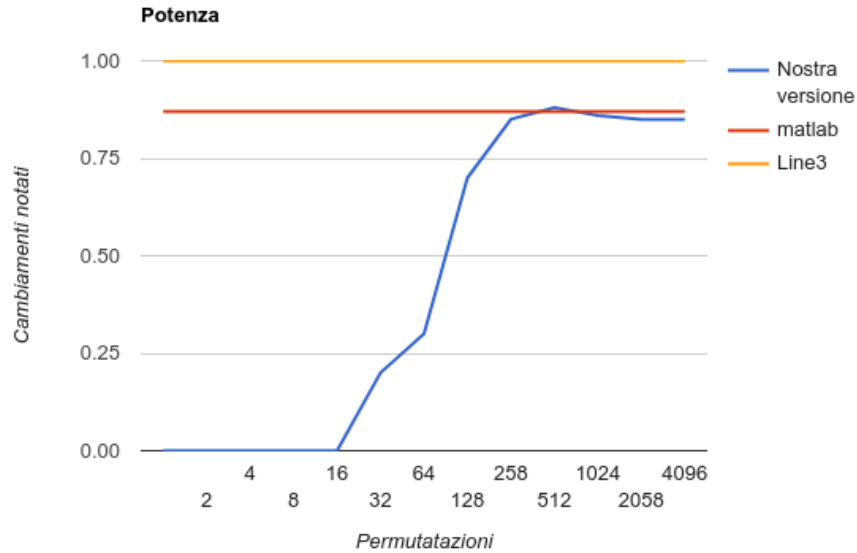
5.1 Correttezza

La correttezza dell'implementazione di un test di ipotesi può essere controllata assicurandosi che in condizioni di ipotesi nulla il test produca esattamente $\alpha\%$ falsi allarmi. Il test deve poi presentare la maggiore potenza possibile in caso di ipotesi contraria. Per utilizzare una metrica di riferimento abbiamo utilizzato i test inclusi nella libreria standard di matlab. Mostriamo quindi il risultati ottenuti dall'uso di una semplice statistica, la differenza delle medie:



La linea rossa è generata da matlab non dipende dal numero di permutazioni, quindi è costante a 0.05, come ci si attende. La linea blu rappresenta invece la nostra versione del permutation test, la quale per bassi valori di permutazione è del tutto incapace di notare cambiamenti nei valori in ingresso. Quando le permutazioni aumentano i falsi allarmi si stabilizzano velocemente ad α .

Ora che siamo certi che sotto ipotesi nulla il nostro algoritmo si comporta correttamente, possiamo quindi analizzare il caso in cui i gruppi provengono da popolazioni diverse.



Come nel caso precedente si nota che pur fallendo per permutazioni in quantità minute risulta che il test produce risultati simili a quello standard. Si noti inoltre che come nel caso precedente non sembra esserci alcun vantaggio nell'utilizzare una grandi numeri di permutazioni, poiché quando raggiungono l'ordine delle centinaia la potenza si stabilizza. Questo significa che, date le ottimizzazioni sull'uso della memoria, se si fissa un numero di permutazioni, è possibile processare set di dati con utilizzi di tempo e memoria che dipendono linearmente solo dalla dimensione dei dati in ingresso. Poiché è evidente che per processare un set di dati è per lo meno necessario valutare almeno una volta tutti gli elementi, ciò risulta essere il migliore risultato possibile.

Ciò che rimane è mostrare la velocità di throughput che è possibile raggiungere è superiore ad un implementazione basata su cpu.

5.2 Lavori Futuri

A seconda di quale sia il campo di applicazione di un test di ipotesi ad alta performance è possibile che ulteriori caratteristiche oltre a quelle indicate siano richieste. Si può immaginare che la necessità di analizzare eventi che avvengono ad alta frequenza sia spesso associata alla necessità di processare tali dati sotto forma di stream continui fino a che non viene notato un cambiamento nella natura dei dati osservati. Una delle soluzioni di tale problema prende il nome di CPM e consiste nell'analizzare i dati una prima volta, per poi spostare il focus dell'attenzione man mano che nuovi dati vengono ricevuti. Il tentativo di integrare efficientemente il permutation test e il cpm si è rivelato fallimentare. Inoltre, anche le integrazioni più semplici sono inefficienti, poiché le ottimizzazioni necessarie ad uno risultano deleterie all'altro. Questo limite pare significativo e ne consegue che la risoluzione possa essere soggetto di ulteriori studi.