

Memoria prácticas - *ALT*

David Arnal García

Luis Alberto Álvarez Zavaleta

Lugman Sami Ahmad Mansilla

Labores realizadas:

En cuanto a las labores realizadas, se ha decidido realizar todas las tareas colaborativamente todos conectados a la vez usando *LiveShare* de *Visual Studio Code*, por lo que estas han sido hechas por todos los miembros del equipo de manera conjunta.

Descripción resumida del trabajo realizado, indicando si se han implementado todos los requerimientos obligatorios:

Para la realización de este proyecto, hemos realizado las partes obligatorias de cada apartado, esto es:

Tarea 1

Hemos hecho la implementación de las tres distancias:

- Distancia de *Levenshtein*
- Distancia de *Damerau-Levenshtein* restringida
- Distancia de *Damerau-Levenshtein* “intermedia” para $cte = 1$

Basándonos en los algoritmos que hemos encontrado en los apuntes de teoría para la distancia de *Levenshtein*. Para la implementación de la distancia restringida e intermedia hemos tomado de base el primer algoritmo de *Levenshtein* y lo hemos modificado para obtener las versiones de *Damerau-Levenshtein* ‘restringida’ e ‘intermedia’.

Tarea 2

Mejoramos las distancias de la tarea 1 añadiendo *thresholds* como una variable que fije un umbral en el cual nuestro algoritmo debe dejar de calcular distancias, si este es alcanzado. Esto lo hemos implementado en *test_tarea2.py*.

Tarea 3

En el fichero *spellsuggest.py* hemos implementado en la función *main*, la estructura de los resultados de los distintos algoritmos de distancia, utilizando el código que se había proporcionado en *plantilla_generar_resultados.py*.

Usamos el comparador de archivos de *Visual Studio* para comprobar que los resultados proporcionados con los ficheros de resultados que se nos habían proporcionado eran los mismos. Además, también implementamos una cota optimista para evitar la ejecución de las distancias de edición de la tarea 2, si esta cota es mayor que el *threshold*.

Tarea 3: cota optimista basada en contar el número de veces que aparece cada letra en cada palabra (opcional)

En el método *suggest_opt()*, dentro de *spellsuggest.py*, hemos implementado esta nueva cota para la distancia de *Levenshtein*, de forma que como se indica en esta en caso de ser alcanzada no se realice la invocación al método de *Levenshtein*.

Tarea 4

Hemos realizado la comparación entre un par de cadenas usando *tries*, en la cual este podrá contener una cadena más corta siendo esta el prefijo de otra más larga, de esta forma, aprovechando la compactación que ofrece tener un árbol para representar las palabras que tendría un diccionario del que deseamos realizar la

comparación con nuestra cadena. En esta parte, recibiendo una cadena y el *trie*, hemos realizado la comparación hacia atrás de *Levenshtein*.

Tarea 4: Versiones hacia atrás para *Damerau-Levenshtein* restringida e intermedia usando *tries* (opcional)

En *test_tarea4.py* hemos implementado tanto *dp_restricted_damerau_trie* como *dp_intermediate_damerau_trie*.

Tarea 5

Hemos realizado un estudio experimental para determinar qué versiones de las anteriores son las más eficientes. Para este estudio, hemos comprobado para unos datos concretos y los distintos algoritmos implementados, la eficiencia de estos con distintas tallas, distintos valores de *threshold*, analizando distintas palabras, todo ello con los mismos datos para poder comparar entre ellos su eficiencia. Comprobaremos la media, mediana y la desviación típica.

En vista de los resultados obtenidos, expuestos en el anexo, hemos decidido usar *Levenshtein* como distancia predeterminada en *SAR_lib*, debido a que es el que mejores resultados de tiempos nos ha proporcionado para todos los *threshold* que hemos usado para su estudio. Adicionalmente, hemos elegido usar *threshold* '3' como opción predeterminada, ya que da un número de errores aceptables para conseguir sugerencias.

Seleccionado unos valores medios que podemos tomar como referencia para apreciar esta elección, podemos apreciar para *threshold* 3 y una talla de 20.000 que podemos considerar como una talla aceptable los siguientes valores:

	levenshtein	levenshtein-intermedia	levenshtein-restricted	levenshtein-trie	levenshtein-intermedia-trie	levenshtein-restricted-trie
Media	1,208	1,494	1,567	2,739	2,925	3,378
Mediana	1,219	1,531	1,594	2,820	2,984	3,477
Desviación típica	0,3412	0,4022	0,4259	0,3877	0,4371	0,5045

Podemos apreciar cómo se obtienen los mejores tiempos con un ajuste más aproximado a los distintos tamaños de las palabras en base a los resultados con *Levenshtein*.

Tarea 5: generar errores ortográficos (opcional)

En el archivo *medirTiempos.py* hemos generado un método llamado *get_random_terms* que toma un número de palabras aleatorias del vocabulario y aplica transformaciones de borrado, intercambio e inserción de palabras para simular errores ortográficos antes de medir tiempos.

Como podemos apreciar, al aumentar la talla y el *threshold*, la implementación de las distancias de *Levenshtein* con *tries* proporciona mejores tiempos.

Tarea 6

Hemos creado un nuevo archivo llamado *SAR_ALT_spellsuggest.py*, siendo una modificación de *spellsuggest.py*, en el que, en lugar de tomar un archivo, toma directamente el vocabulario. Con esto, podremos guardar en *self.spellsuggester* el vocabulario de cada campo, el cual se guardará en el *bin* generado por el *Indexer* y podrá ser llamado para ser utilizado desde el *Searcher*.

Para ejecutarlo, podremos usar los métodos usados en SAR:

```
py .\SAR_Indexer.py -S -P -M -O -U corpora/2015 2015_index_full.bin
```

donde *U* es la extensión para incluir el *suggester* en el indexador.

```
py .\SAR_Searcher.py -C -Q 'word' -E -D 'levenshtein' -H 'threshold'  
2015_index_full.bin
```

donde *-E* activa la búsqueda aproximada, *-D* la distancia y *-H* el *threshold*, por defecto utiliza *Levenshtein* con *threshold = 3*.

Una pequeña "demo" del funcionamiento del indexador y del buscador. Se pueden incluir capturas de pantalla:

```
py ./exec.py
```

Hemos realizado un pequeño *script* “*exec.py*” que realiza la ejecución de los distintos programas, de forma que ejecuta *SAR_indexer.py* para la creación del archivo *.bin*, que es el que emplea *SAR_Searcher.py* y después ejecuta este, indicándoles las distintas palabras a comprobar y los distintos *threshold* a probar.

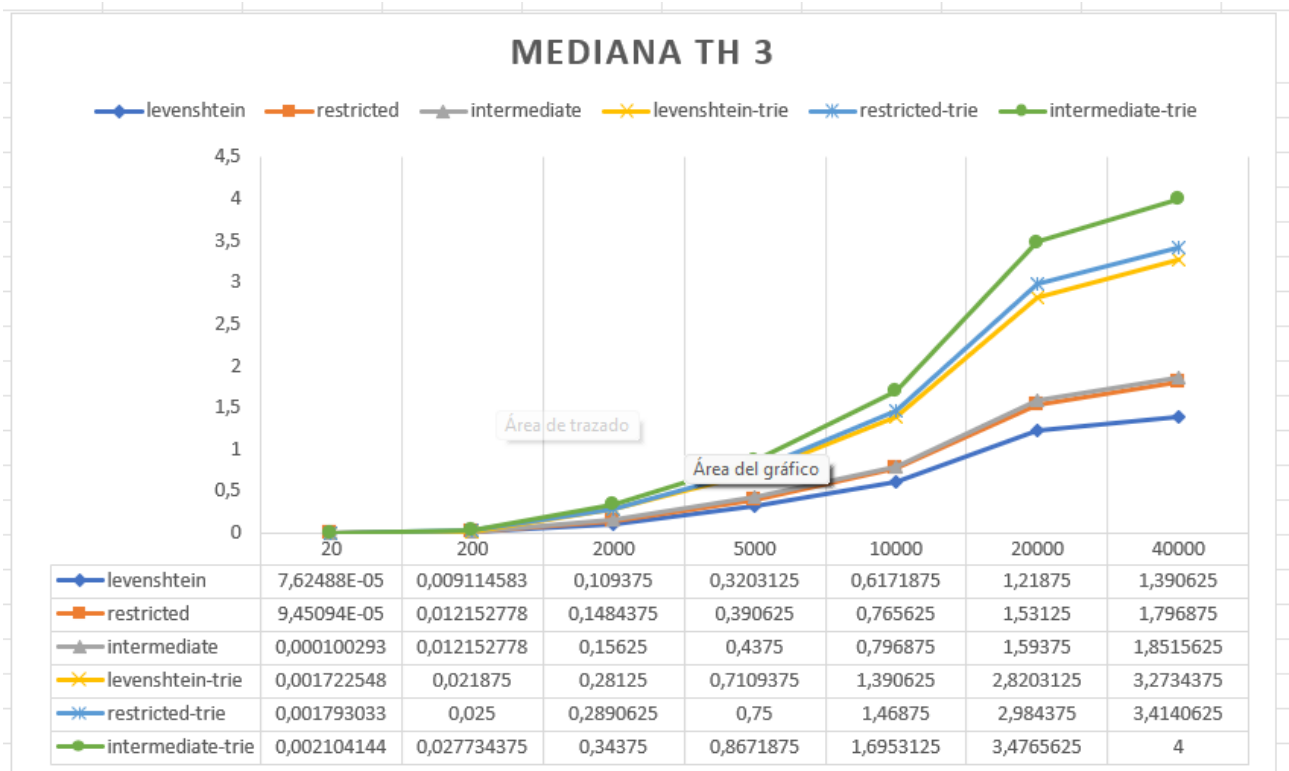
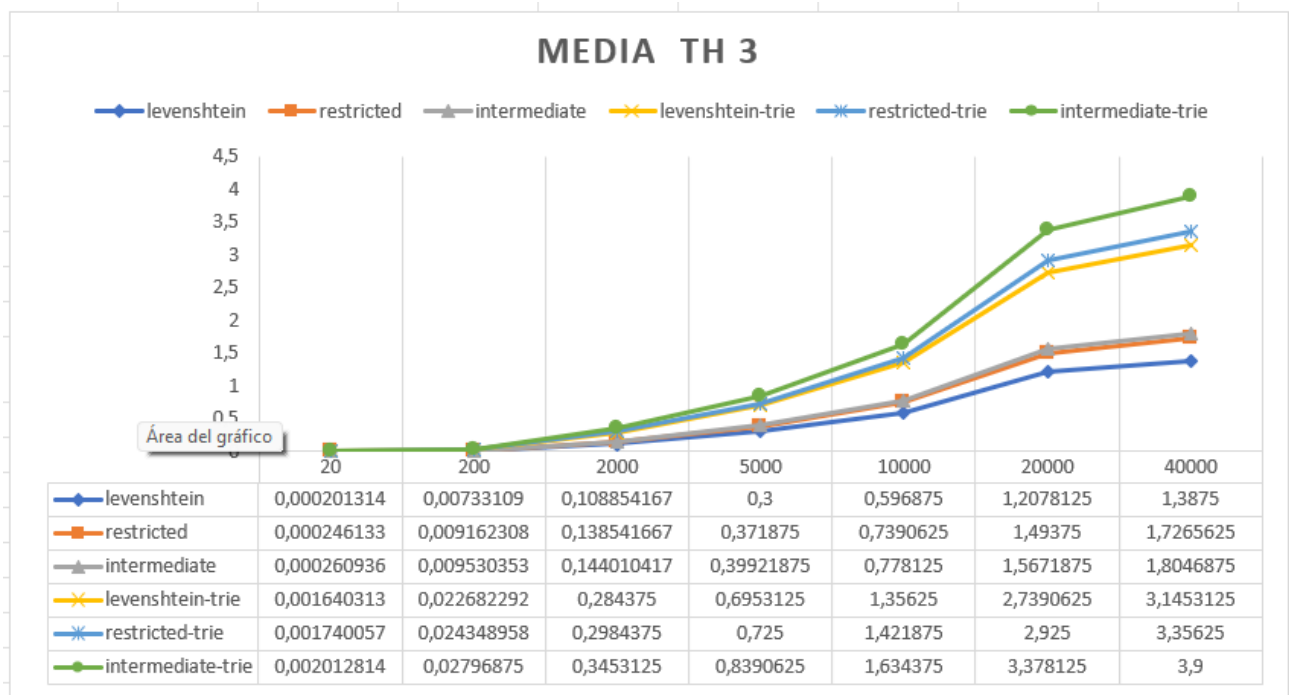
Aquí podemos apreciar una imagen de nuestro *script*:

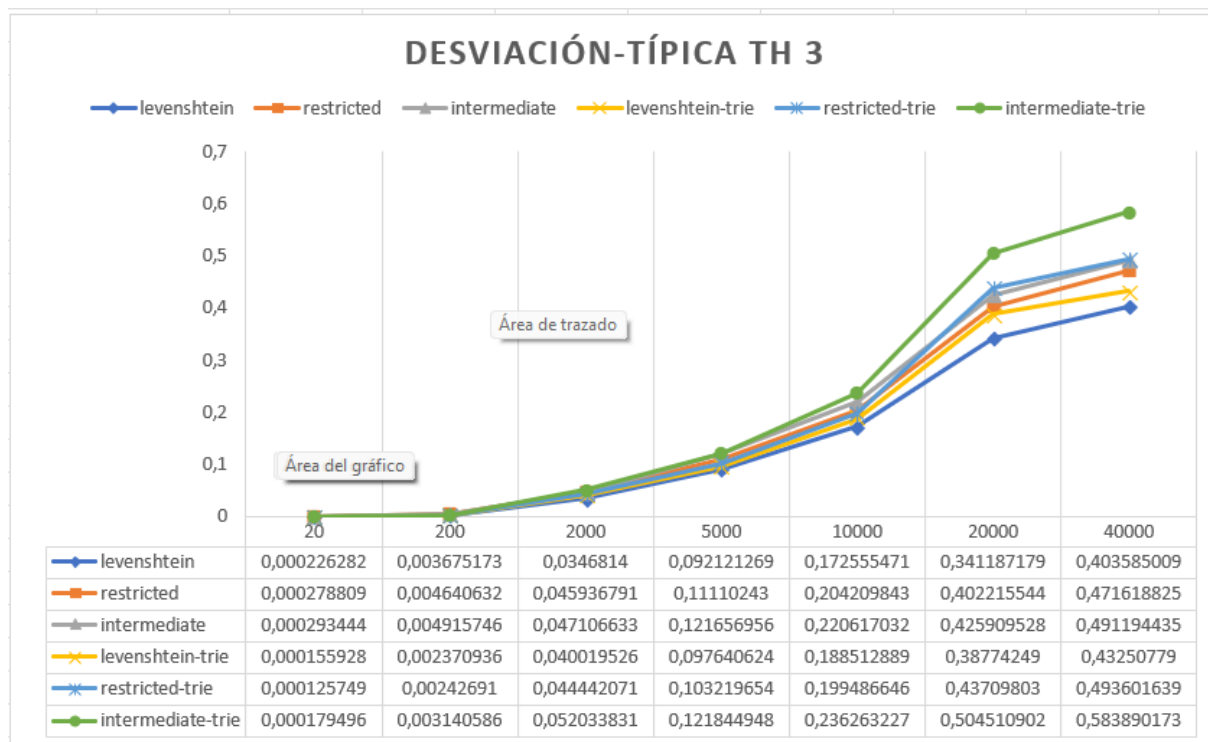
```
import os
os.system('py .\SAR_Indexer.py -S -P -M -O -U corpora/2015 2015_index_full.bin')
for i in range(1, 4):
    for dis in ["levenshtein"]:
        print("TH: %d" % i)
        print(dis)
        for word in ["alexanderx", "iluminaciónx", "capitla", "secuelax",
                    "asoman2", "virneg", "perwisos", "persimos", "limertad", "sagardo",
                    "aperturaq", "loursed", "cuagro", "planrta", "NOT antihéroes AND ecuador",
                    "blindarse OR árboles AND perjuicio", "orígenes AND aversión",
                    "actualmente AND NOT reubicar", "ventanasse OR rebuscar",
                    "tortuosa OR atardecer", "NOT alexander AND iluminación",
                    "capital OR secuela AND asoman"]:
            os.system('py .\SAR_Searcher.py -C -Q "' + word + ' ' +
                    '" -E -D ' + dis + ' -H ' + str(i) + ' 2015_index_full.bin')
```

Obteniendo, así, los resultados de la ejecución de la búsqueda aproximada en nuestro buscador de noticias, que serán los que compararemos con los esperados para estos algoritmos y con los *threshold* indicados, viendo cómo coinciden con los proporcionados en *PoliformaT*. A continuación, mostraremos una imagen de la salida de este para una pequeña cantidad de valores:

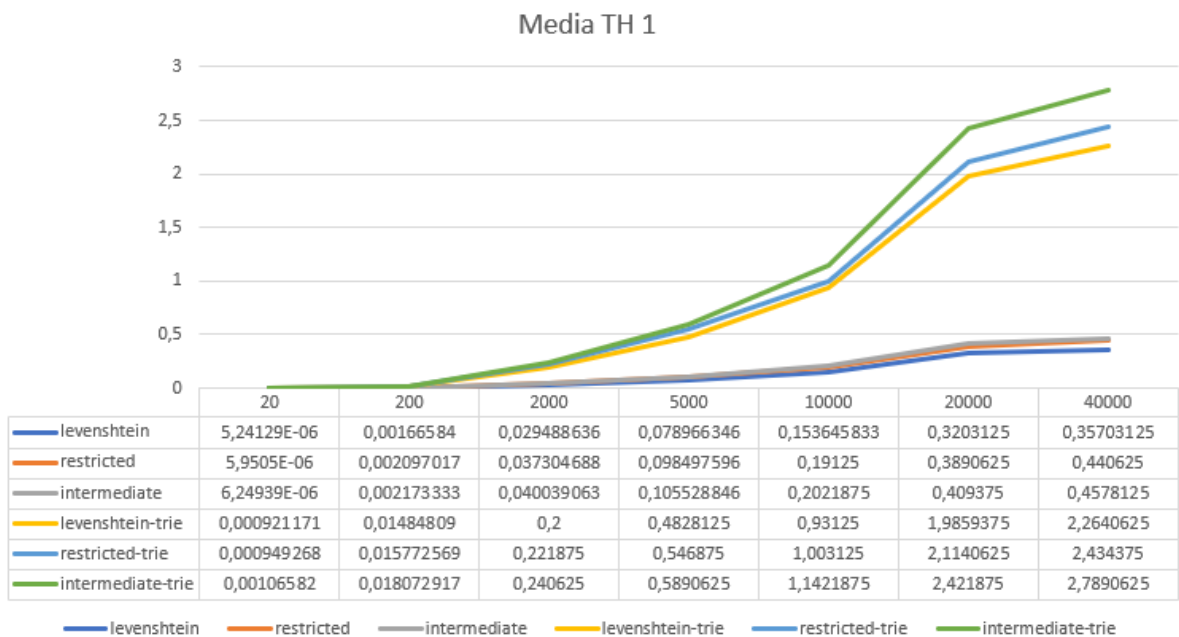
```
TH: 1
levenshtein
alexanderx      2
iluminaciónx    8
capitla 10
secuelax        2
asoman2 2
virneg 0
perwisos        6
persimos        1
limertad        58
sagardo 0
aperturaq       26
loursed 0
cuagro 184
planrta 50
NOT antihéroes AND ecuador      11
blindarse OR árboles AND perjuicio 0
orígenes AND aversión 1
actualmente AND NOT reubicar    73
ventanasse OR rebuscar 2
tortuosa OR atardecer 9
NOT alexander AND iluminación 8
capital OR secuela AND asoman 0
TH: 2
```

Media, mediana y desviación típica para un *threshold* de 3

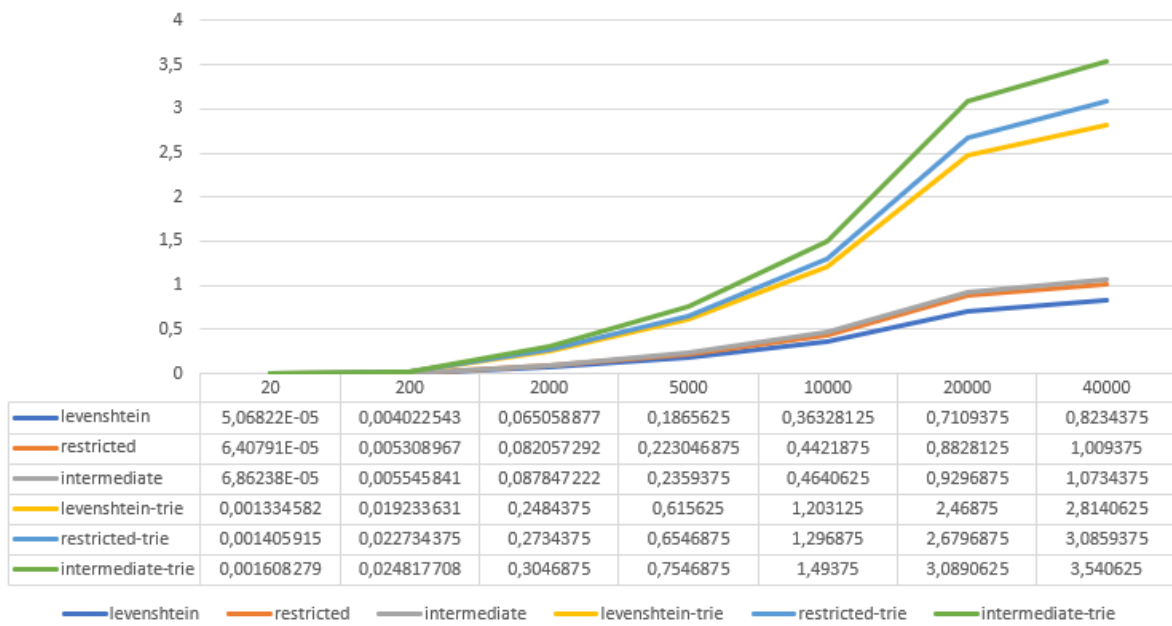




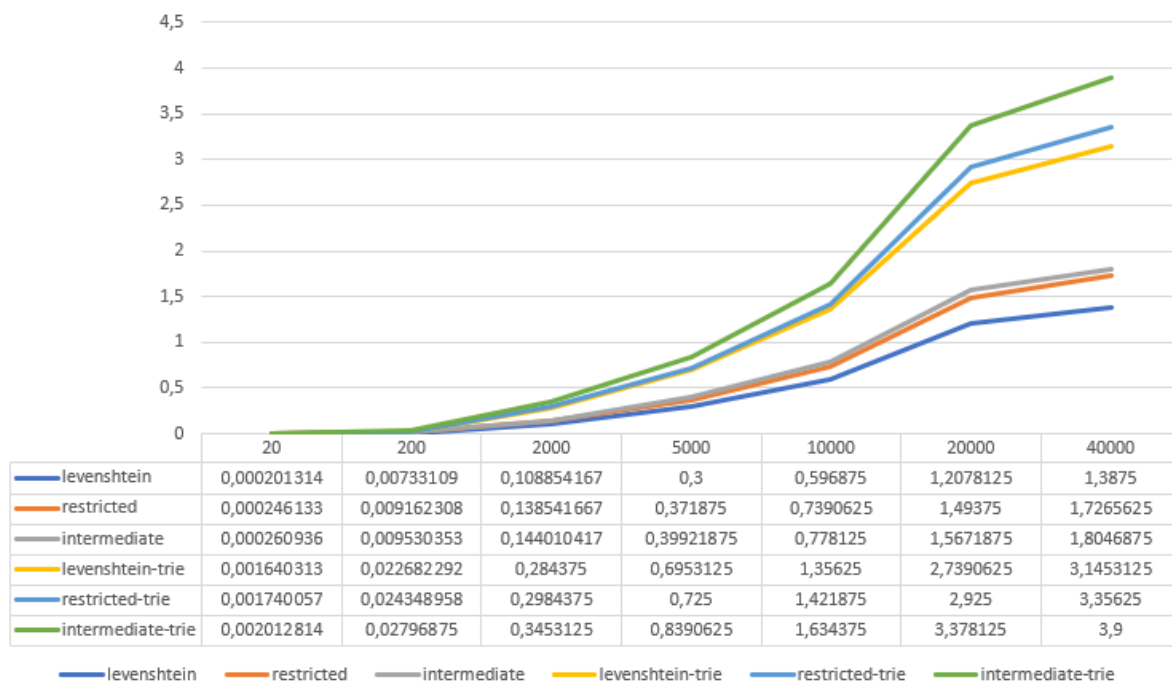
Media de los distintos algoritmos



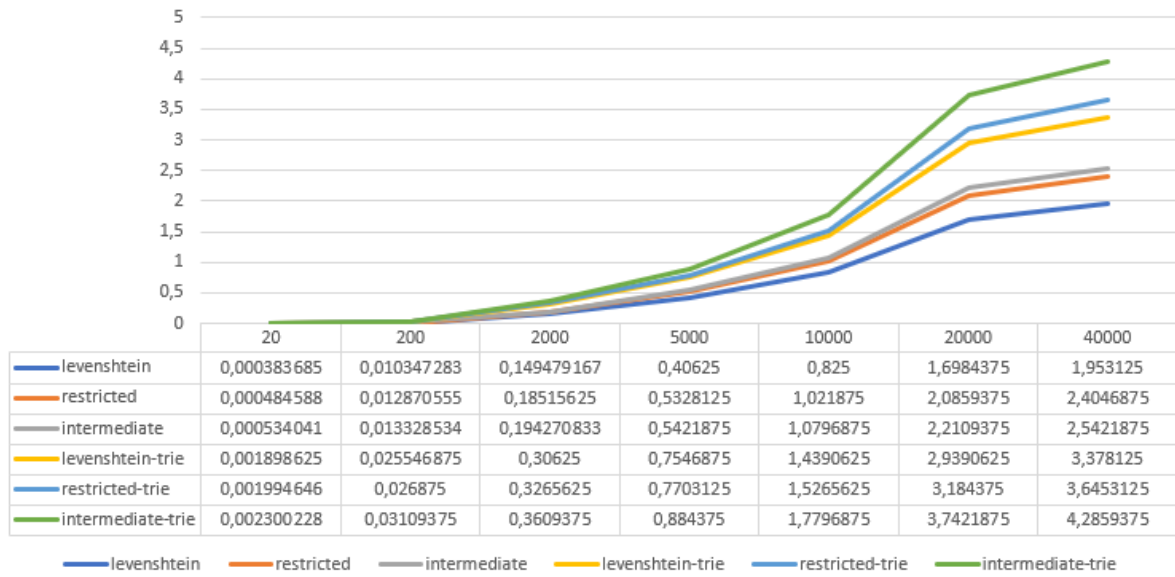
Media TH 2



Media TH 3



Media TH 4



Media TH 5

