

Universidad de Murcia
INGENIERÍA INFORMÁTICA



TECNOLOGÍAS ESPECÍFICAS DE LA INGENIERÍA INFORMÁTICA
Boletín de prácticas de los bloques 3 y 4

Autor: Victorio J. Molina Bermejo
DNI: 48632380-F Grupo: 3.2

FACULTAD DE INFORMÁTICA

Junio 2020

Índice

| | |
|---|--------------|
| 1. Especificación de la práctica | III |
| 2. Resolución de la práctica | IV |
| 2.1. Código | IV |
| 2.1.1. Código Python | IV |
| 2.1.2. Código C | XI |
| 3. Guía para la ejecución del programa | XV |
| 4. Herramientas utilizadas durante el desarrollo | XV |
| 4.0.1. Numpy | XV |
| 4.0.2. Matplotlib | XVI |
| 4.0.3. Ctypes | XVI |
| 4.0.4. Profiling | XVII |
| 4.0.5. Git | XVII |
| 4.0.6. GitHub | XVIII |
| 5. Análisis del orden de complejidad de la función C | XVIII |
| 6. Resultados del profiling | XIX |
| 6.1. Memory profiling | XIX |
| 6.1.1. Linear profiling y gráficas generadas con Matplotlib . . . | XXI |
| 7. Bibliografía | CXXII |

1. Especificación de la práctica

Imaginemos que dada, una lista de elementos, queremos obtener una nueva versión de la misma sin repeticiones. Python nos ofrece los conjuntos para ello, pero también podríamos conseguir lo mismo usando listas, diccionarios, etc. Queremos evaluar las distintas posibilidades que tenemos en Python para saber cuál es la más eficiente y compararlas con una implementación propia en C de una función de eliminación de duplicados. Para ello, debemos implementar un programa de Python que reciba tres parámetros: un fichero de entrada con una lista de elementos, uno por línea, un fichero de salida donde se guardará esa misma lista de elementos, pero sin repeticiones, y un fichero de salida donde se guardará en PDF la gráfica que mostrará el tiempo usado por cada técnica para distintos tamaños de datos; esta gráfica se generará usando Matplotlib. Un detalle a tener en cuenta es que el orden de los elementos en el fichero de salida puede ser distinto de aquel que tienen en el fichero de entrada.

Podemos suponer que los elementos del fichero de entrada son números naturales entre 0 y un cierto límite. Para hacer las pruebas, se recomienda generar un fichero que contenga 200000 números naturales aleatorios entre 0 y 99999 (ambos inclusive) y hacer un estudio de cada técnica con los primeros 2000 elementos del fichero, los primeros 4000 elementos del fichero, los primeros 6000 elementos del fichero, y así sucesivamente hasta el total de 200000 elementos. Para generar el fichero de entrada, se sugiere hacer un pequeño programa en Python que reciba como parámetro el nombre del fichero y, usando Numpy, se generen los 200000 números naturales entre 0 y 99999 (el número de elementos a generar también podría ser un parámetro de entrada).

Al menos, hay que analizar dos formas de obtener una lista sin repeticiones usando dos estructuras de datos distintas de Python y una tercera forma que debemos implementar en código en C y que usaremos desde Python a través de ctypes. En esta tercera forma, la función de C a implementar debe recibir la lista de elementos original y devolver la lista de elementos sin duplicados. La memoria dinámica que sea necesaria reservar para el intercambio de información entre el código de Python y el de C se reservará desde Python. Los alumnos tienen que decidir qué código de C implementar para obtener el resultado deseado de la forma más eficiente posible.

Además de la gráfica generada en PDF, que habrá que incluir en el documento LATEX final donde se explique el trabajo realizado, también habrá que incluir el resultado del perfilado (profiling) del tiempo de ejecución de cada técnica de eliminación de duplicados usada y una gráfica que muestre el consumo de memoria por parte del programa durante su ejecución. Para realizar estos perfilados, pueden usarse los módulos *line_profiler* y *memory_profiler* descritos en teoría. En el perfilado del tiempo de ejecución no es necesario hacer un estudio para cada posible tamaño de la lista de elementos de la que hay que eliminar duplicados. Basta con que el perfilado se haga de tal manera que refleje lo mismo que la gráfica generada en PDF en cuanto al orden de las distintas técnicas en lo que a eficiencia se refiere.

2. Resolución de la práctica

2.1. Código

2.1.1. Código Python

A continuación se presenta el código Python utilizado para resolver el enunciado, desde el módulo main, y el fichero contenedor de las funciones de eliminación de duplicados, en la que se incluye el Wrapper que permite usar la biblioteca C existente, hasta el script para la generación del fichero de entrada.

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  '''
5  Created on Thu May  7 01:40:07 2020
6
7  @author: VictorioMolina
8  '''
9
10 import sys
11 import utils
12 import matplotlib.pyplot as plt
13 import numpy as np
14 import time
15 from inspect import getmembers, isfunction
16 from matplotlib.backends.backend_pdf import PdfPages
17 from matplotlib import gridspec
```

```

18
19 plt.rcParams.update({'font.size': 14})
20 plt.rcParams.update({'figure.max_open_warning': 0})
21
22
23 # Method that turns the file lines into a list
24 def file_to_list(path):
25     try:
26         infile = open(path, 'r')
27     except IOError:
28         # The file doesn't exist
29         print("The file {} doesn't exists".format(path))
30         sys.exit(-1)
31     else:
32         # The file exists
33         elements = []
34
35         # Add the file lines to the list of elements
36         elements.extend(int(line) for line in infile)
37
38         infile.close()
39         return elements
40
41 # Method that writes the result of the program in an specific
42 # file
43 def result_to_file(result, path):
44     try:
45         outfile = open(path, 'x')
46     except FileExistsError:
47         print("The file {} already exists".format(sys.argv[2]))
48         sys.exit(-1)
49     else:
50         for element in result:
51             outfile.write(str(element) + '\n')
52         outfile.close()
53
54 # Method that generates a plot and save it in an specific PDF
55 # file
56 def generate_plots(title, x, y, pdf):
57     fig = plt.figure(figsize=(14, 19.8)) # A4 Size
58
59     fig.suptitle(title, fontsize=30)
60
61     spec = gridspec.GridSpec(ncols=3, nrows=2, width_ratios
62                             =[0.5, 4, 0.5])

```

```

60
61 # Bar plot
62 ax = fig.add_subplot(spec[1])
63 plt.xlabel("Execution Time")
64 width = 0.75 # Bars' width
65 ind = np.arange(len(y))
66 ax.barh(ind, y, width, color="blue")
67 ax.set_yticks(ind + width / 4)
68 ax.set_yticklabels(x, minor=False)
69
70 # Pie chart
71 ax1 = fig.add_subplot(spec[4])
72 explode = (0, 0, 0, 0, 0, 0.05) # only "explode" the 6th
73 slice which represents the C function
74 ax1.pie(
75     [_ * 1000 for _ in y],
76     explode=explode,
77     labels=x,
78     autopct='%1.1f%%',
79     shadow=False,
80     startangle=90
81 )
82 ax1.axis('equal') # Equal aspect ratio ensures that pie is
83 drawn as a circle.
84
85 # Save current figures in a pdf page
86 pdf.savefig()
87
88 # Method that studies the efficiency of each technique for
89 different data sizes
90 # generating graphics that will be saved in a pdf
91 def study_techniques_efficiency(seq, path):
92     try:
93         pdf_outfile = open(path, 'x')
94     except FileExistsError:
95         print("The file {} already exists".format(path))
96         sys.exit(-1)
97     else:
98         # Create a multipage PDF file
99         pdf = PdfPages(path)
100
101         # Get all the functions of the module
102         functions_list = [o for o in getmembers(utils,

```

```

102         isfunction)]
103         # Make an study of each technique in size intervals of
2000 elements
104         size_interval = 2000
105         for i in range(int(len(seq) / size_interval)):
106             # For each technique, calling its function with the
corresponding time taking
107             cutten_list = seq[0:(i+1)*size_interval]
108             exec_times = [] # Will store the exec time of each
function
109
110             for function_name, function_obj in functions_list:
111                 # Calculate the exec time
112                 t0 = time.time_ns()
113                 result = function_obj(cutten_list)
114                 t_exec = (time.time_ns() - t0) / 1.0e9
115
116                 # Store the ecex time
117                 exec_times.append(t_exec)
118                 print("{} has taken {} seconds to execute".
format(function_name, t_exec))
119
120             # Generate a plot with the execution time stats
121             generate_plots(
122                 "TIMES FOR THE FIRST {} NUMBERS".format((i
+1)*size_interval),
123                 [f[0] for f in functions_list],
124                 exec_times,
125                 pdf
126             )
127
128             # Close the pdf
129             pdf.close()
130
131 def main():
132     """
133
-----
134     PROGRAM ARGUMENTS:
135     -> infile: input file with a list of elements, one
per line
136     -> outfile: output file where that same list of
elements will be saved,

```

```

137         but without repetitions
138         -> pdf_outfile: output file where the graph that
            will show the time used
139         by each technique for different data sizes
            will be saved in pdf
140
141         -----
142
143         """
144         # Command line argument control
145         if len(sys.argv) != 4:
146             print("Usage: {} infile outfile pdf_outfile".format(sys.
argv[0]))
147             sys.exit(0)
148
149         # Reading the sequence of the input file
150         seq = file_to_list(sys.argv[1])
151
152         # Calling the external function through its wrapper
153         result = utils.remove_duplicates_6(seq)
154
155         # Writing the result in the outfile
156         result_to_file(result, sys.argv[2])
157
158         # Finally, study the different techniques efficiency
159         study_techniques_efficiency(seq, sys.argv[3])
160
161 if __name__ == '__main__':
162     main()

```

File 1: main.py


```

1 # -*- coding: utf-8 -*-
2
3 '''
4 This file contains a number of useful functions for removing
5     duplicates.
6
7 @author: VictorioMolina
8 '''
9
10 import ctypes, os
11 from functools import reduce
12 import numpy as np
13 import time
14
15 # Loading the shared library into ctypes
16 LIBREMOVE_DUPLICATES = ctypes.CDLL(
17     os.path.abspath(
18         os.path.join(os.path.dirname(__file__), "../C/
19             libremove_duplicates.so.1")
20     )
21 )
22
23 # @profile
24 def remove_duplicates_1(seq):
25     # Not preserving the order
26     return list(set(seq))
27
28 # @profile
29 def remove_duplicates_2(seq):
30     # Preserving the order
31     found = set()
32     return [x for x in seq if not (x in found or found.add(x))]
33
34 # @profile
35 def remove_duplicates_3(seq):
36     # Not preserving the order
37     keys = {}
38     for x in seq:
39         keys[x] = 1
40     return list(keys.keys())
41
42 # @profile
43 def remove_duplicates_4(seq):
44     # Not preserving the order and using NumPy

```

```

44     return list(np.unique(seq))
45
46 # @profile
47 def remove_duplicates_5(seq):
48     # Preserving the order and using NumPy
49     indexes = sorted(np.unique(seq, return_index=True)[1])
50     return [seq[i] for i in indexes]
51
52 # Python Wrapper for calling the C function
53 # @profile
54 def remove_duplicates_6(seq):
55     # Object corresponding to the function within the library
56     func_remove_duplicates = LIBREMOVE_DUPLICATES.
57     remove_duplicates
58
59     # Function prototype
60     func_remove_duplicates.argtypes = [
61         ctypes.POINTER(ctypes.c_uint),
62         ctypes.POINTER(ctypes.c_uint),
63         ctypes.c_uint
64     ]
65
66     # This function returns void
67     func_remove_duplicates.restype = None
68
69     # Variable in which we will collect the result of the
70     # function.
71     result = (ctypes.c_uint * len(seq))()
72
73     # Call to the shared library function.
74     func_remove_duplicates((ctypes.c_uint * len(seq))(*seq),
75                             result, len(seq))
76
77     # Copying result of the function to another vector
78     vout=[*result] # This type of copy is the 'fastest' in
79     # Python
80
81     # Removing all zeros least the first one
82     while len(vout) > 1 and vout[-1] == 0:
83         vout.pop()
84
85     # Finally, we return the output vector
86     return vout

```

File 2: /src/Python/utils.py

```

1 # Generate a file made up of random integers
2
3 import numpy as np
4
5 quantity = int(input("How many? "))
6 limit = int(input("Max integer: "))
7 filename = "numbers.txt"
8
9 numbers = np.random.randint(low=0, high=limit+1, size=quantity)
10
11 f = open(filename, 'w')
12 for n in numbers:
13     f.write("{}\n".format(n))

```

File 3: /test/generate_test_file.py

2.1.2. Código C

En lo que al código C respecta, he implementado un fichero de cabecera con los prototipos de función, un fichero .c con la implementación del método de eliminación de duplicados dando uso del algoritmo *Heap Sort*, además de generar un archivo Makefile para la compilación del código con la biblioteca gcc, que conlleva a la creación de la biblioteca compartida nativa de Linux *libremove_duplicates.so.1.0.1* utilizada en el wrapper de Python.

```

1 #ifndef __REMOVE_DUPLICATES_H__
2 #define __REMOVE_DUPLICATES_H__
3
4 /* Working with memory already reserved by Python! */
5
6 // Function Prototypes
7
8 void remove_duplicates(const unsigned int *vin,
9                       unsigned int *vout,
10                      const unsigned int size);
11
12 void heap_sort(unsigned int arr[], const unsigned int size);
13
14 void create_max_heap(unsigned int arr[],
15                     const unsigned int size,
16                     unsigned int root_index);
17
18 void swap_nodes(unsigned int *ptr1, unsigned int *ptr2);

```

```
19
20 #endif
```

File 4: /src/C/remove_duplicates.h

```
1 #include "remove_duplicates.h"
2
3 // Function to remove duplicated elements in an array
4 void remove_duplicates(const unsigned int *vin, unsigned int *
   vout, const unsigned int size)
5 {
6     /*
7
8         Complexity Analysis
9         -----
10        TIME                 $O(n * \log(n)) + O(n) = O(n * \log(n))$ 
11        SPACE                 $O(n) + O(1) = O(n)$ 
12        -----
13
14    */
15
16    // Creating copy of the entry vector... (good choice)
17    unsigned int arr[size];
18    for(int i = 0; i < size; i++)
19    {
20        arr[i] = vin[i];
21    }
22
23    // Sorting the array with the heap-sort algorithm
24    heap_sort(arr, size);
25
26    int i = 0;
27    int j = 0;
28    while (i < size)
29    {
30        // Insert element to the result vector
31        vout[j] = arr[i];
32        j = j + 1;
33
34        while (arr[i] == arr[i + 1])
35        {
36            i = i + 1;
37        }
38        i = i + 1;
```

```

39     }
40 }
41
42 // Heap-sort main function
43 void heap_sort(unsigned int arr[], const unsigned int size)
44 {
45     /*
46
47         Complexity Analysis
48         -----
49         best, worst, average TIME      O(n * log(n))
50         SPACE                          O(1)
51         -----
52
53     */
54
55     for (signed int i = size / 2 - 1; i >= 0; i--)
56     {
57         create_max_heap(arr, size, i);
58     }
59
60     for (signed int i = size - 1; i > 0; i--)
61     {
62         swap_nodes(&arr[0], &arr[i]);
63         create_max_heap(arr, i, 0);
64     }
65 }
66
67 // Heap-sort helper function
68 void create_max_heap(unsigned int arr[], const unsigned int size
69 , unsigned int root_index)
70 {
71     unsigned int root = root_index;
72     unsigned int left = 2 * root_index + 1;
73     unsigned int right = 2 * root_index + 2;
74
75     // If the left son is greater than the root
76     if (left < size && arr[left] > arr[root])
77     {
78         root = left;
79     }
80
81     // If the right son is greater than the root
82     if (right < size && arr[right] > arr[root])
83     {

```

```

83     root = right;
84 }
85
86 // If the root has changed
87 if (root != root_index)
88 {
89     swap_nodes(&arr[root_index], &arr[root]);
90
91     // Continue resolving the new tree
92     create_max_heap(arr, size, root);
93 }
94 }
95
96 // Heap-sort helper function
97 void swap_nodes(unsigned int *ptr1, unsigned int *ptr2)
98 {
99     unsigned int tmp = *ptr1;
100    *ptr1 = *ptr2;
101    *ptr2 = tmp;
102 }

```

File 5: /src/C/remove_duplicates.c

```

1 all: libremove_duplicates.so.1.0.1 libremove_duplicates.so.1
2
3 libremove_duplicates.so.1.0.1: remove_duplicates.c
4     remove_duplicates.h
5     gcc -Wall -c -fPIC -g remove_duplicates.c
6     gcc -shared -Wl,-soname,libremove_duplicates.so.1 -o
7     libremove_duplicates.so.1.0.1 remove_duplicates.o
8
9 libremove_duplicates.so.1: libremove_duplicates.so.1.0.1
10    ln -s libremove_duplicates.so.1.0.1 libremove_duplicates.so.1
11
12 clean:
13    rm -fr libremove_duplicates.so.1.0.1 libremove_duplicates.so.1
14    remove_duplicates.o *~

```

File 6: /src/C/Makefile

3. Guía para la ejecución del programa

Para la ejecución del programa basta con teclear siguiente comando en el terminal:

```
python3 main.py infile outfile pdf_outfile
```

Siendo *infile* el fichero de entrada que contiene el listado de números; *outfile* la ruta del fichero (inexistente) en el que se guardará el resultado del programa, es decir, la lista de numeros del fichero de entrada pero sin duplicados; y *pdf_outfile* la ruta del fichero (también inexistente) en el que se guardarán las gráficas estadísticas creadas durante el estudio de los tiempo de ejecución de las distintas técnicas de eliminación de duplicados implementadas en el programa.

Cabe destacar que he usado en el módulo main he usado `#!/usr/bin/env python3` para la portabilidad en diferentes sistemas en caso de que tengan el intérprete de idiomas instalado en diferentes ubicaciones, lo cual además nos permite ejecutar el programa de la siguiente manera:

```
./main.py infile outfile pdf_outfile
```

4. Herramientas utilizadas durante el desarrollo

4.0.1. Numpy

Numpy es el paquete fundamental para la computación científica con Python. Entre otras cosas contiene:

- Un poderoso objeto de matriz N-dimensional
- Funciones sofisticadas de broadcasting, que permite trabajar con arrays de diferentes tamaños
- Herramientas para integrar código C/C++
- Álgebra lineal útil, transformación de Fourier, y capacidades para la generación de números aleatorios

Además, una gran cantidad de paquetes de cómputo científico se apoyan de Numpy como biblioteca básica, siendo Matplotlib, biblioteca de la que hablaré en la siguiente sección, una de ellas.

En mi caso, me ha sido de utilidad en la eliminación de duplicados, apoyandome en las técnicas que proporciona Numpy, y la generación de numeros aleatorios.

4.0.2. Matplotlib

Matplotlib es el estándar de facto para realizar todo tipo de gráficos en Python, y por tanto resulta obvio que, dada la capacidad de configuración y la calidad de resultados que esta nos ofrece, durante el estudio de los tiempos de ejecución, para la generación de las gráficas para los distintos tamaños de datos, la haya utilizado.

He optado por "plotear" la información en gráficos de barras, ya que es la mejor opción para entender visualmente, de manera clara y sencilla, la comparativa de tiempos, y de tarta, para representar el porcentaje del tiempo total que utiliza cada técnica.

También he importado *PdfPages* de *matplotlib.backends.backend_pdf* para crear el pdf en el que guardar todos los gráficos obtenidos.

4.0.3. Ctypes

Ctypes permite usar bibliotecas existentes (programadas con otros lenguajes) escribiendo simples wrappers de código en Python. Básicamente sirve como puente o interfaz entre el código Python y el código C, pues nos proporciona el mecanismo para realizar la llamada a la función C, mediante sus propios tipos de datos.

Antes de realizar la llamada debemos indicar el prototipo de la función externa (a la que se va invocar) mediante los atributos *argtypes* y *restype*.

4.0.4. Profiling

En ingeniería de software el análisis de rendimiento, comúnmente llamado profiling, es la investigación del comportamiento de un programa usando información reunida desde el análisis dinámico del mismo. Esta técnica nos permite conocer el consumo de memoria de un programa, frecuencia y duración de las llamadas a funciones... ayudándonos a optimizarlo en algunos aspectos.

El profiling del tiempo de ejecución lo he realizado con el módulo *line_profiler*, mientras que el de la memoria lo he hecho con *memory_profiler*.

En la jerarquía de directorios del programa, concretamente en el directorio */doc/profiling* encontramos los distintos ficheros (visualizables) *main.py.lprof* y *mprofile_20200519085312.dat* que contienen el resultado del análisis del tiempo y memoria utilizada, respectivamente.

4.0.5. Git

Como muchas de las grandes cosas en esta vida, Git comenzó con un poco de destrucción creativa y una gran polémica.

En el 2005, la relación entre la comunidad que desarrollaba el kernel de Linux y la compañía que desarrollaba BitKeeper se vino abajo y la herramienta dejó de ser ofrecida de manera gratuita. Esto impulsó a la comunidad de desarrollo de Linux (y en particular a Linus Torvalds, el creador de Linux) a desarrollar su propia herramienta basada en algunas de las lecciones que aprendieron mientras usaban BitKeeper.

Para el control de versiones del proyecto he utilizado esta herramienta, pues me ha permitido manipular y gestionar el código en la línea de tiempo, reflejando los distintos cambios realizados. Además cuenta con una versión gráfica llamada gitg, en la que he podido ir apreciando con los ojos los estados de las ramas que he ido creando.

En la siguiente figura muestro el estado final del proyecto, justo antes de importar este propio fichero Latex en la carpeta de documentación.

| Refactoring | Victorio Molina | Hace 3 horas |
|--|-----------------|---------------|
| Merge branch 'develop' | Victorio Molina | Hace 4 horas |
| Tip for profiling memory consumption | Victorio Molina | Hace 4 horas |
| Refactoring project | Victorio Molina | Hace 4 horas |
| Time profiling, memory profiling, and final results. | Victorio Molina | Hace 4 horas |
| Method which studies all techniques efficiency generating plots into a PDF | Victorio Molina | Hace 5 horas |
| Reading the input file and checking its existence | Victorio Molina | Hace 11 horas |
| Script for the test file generation | Victorio Molina | Hace 11 horas |
| Update utils.remove_duplicates_8 | Victorio Molina | 8 may, 00:21 |
| I forgot to use the sorted array for the duplicates removal... FIXED! | Victorio Molina | 7 may, 22:17 |
| I forgot to use the sorted array for the duplicates removal... FIXED! | Victorio Molina | 7 may, 22:17 |
| Solved bug in the heap sort function | Victorio Molina | 7 may, 21:11 |
| Merge branch 'develop' into fix | Victorio Molina | 7 may, 21:05 |
| Update utils.py | Victorio Molina | 7 may, 21:04 |
| Implementation of the Python Wrapper | Victorio Molina | 7 may, 19:31 |
| Update .gitignore | Victorio Molina | 7 may, 18:39 |
| Revert "Revert 'Another hotfix'" | Victorio Molina | 7 may, 18:36 |
| Revert "Another hotfix" | Victorio Molina | 7 may, 18:34 |
| Another hotfix | Victorio Molina | 7 may, 18:27 |
| Another hotfix | Victorio Molina | 7 may, 18:27 |
| Fixed remove_duplicates.c | Victorio Molina | 7 may, 18:20 |
| Fixed remove_duplicates.c | Victorio Molina | 7 may, 18:20 |
| C Makefile | Victorio Molina | 7 may, 17:53 |
| C duplicates removal function implemented with a Heap Sort | Victorio Molina | 7 may, 07:25 |
| Added C header file | Victorio Molina | 7 may, 05:19 |
| Update project structure and .gitignore | Victorio Molina | 7 may, 04:04 |
| I have implemented a few functions in python for deleting duplicates in an 'utils' module, and refactored the project structure to match it to the practice specification. | Victorio Molina | 7 may, 03:52 |
| Preparation of the main file. Control of arguments received by command line. | Victorio Molina | 7 may, 02:22 |
| First definition of 'README.md' and '.gitignore' | Victorio Molina | 7 may, 01:54 |

Figura 1: Estado final del proyecto git

4.0.6. GitHub

He subido el proyecto a mi repositorio online de GitHub. Este es el enlace al mismo: <https://github.com/VictorioMolina/Python-remove-duplicates-profiling>

5. Análisis del orden de complejidad de la función C

El algoritmo de eliminación de duplicados que he implementado es sencillo, primero ordenamos el vector de entrada (para así poder eliminar los duplicados de una pasada) mediante uno de los mejores algoritmos de ordenación en lo que a tiempo y memoria se refiere, para luego, como ya he dicho, realizar la tarea de unificar la estructura.

El algoritmo de ordenación que he utilizado es el *Heap Sort*, el cual tiene un orden de complejidad, **para todos los casos**, de $n \cdot \log(n)$ refiriendonos al tiempo de ejecución, y constante en cuanto a memoria. Dentro de lo que cabe es un buen orden, pues obtenemos más eficiencia para tamaños de datos gigantescos en comparación al resto de algoritmos.

Sumando el peso de la ordenación y de la unificación de la estructura, finalmente obtenemos un orden de complejidad $n \cdot \log(n)$ para el tiempo, y n para la

memoria. Este último podría haber sido constante si no hubiese realizado la copia del vector que recibimos como argumento, pero he preferido no realizarla para evitar la destrucción del vector original y realizar los ajustes necesarios para que Python no se vuelva loco.

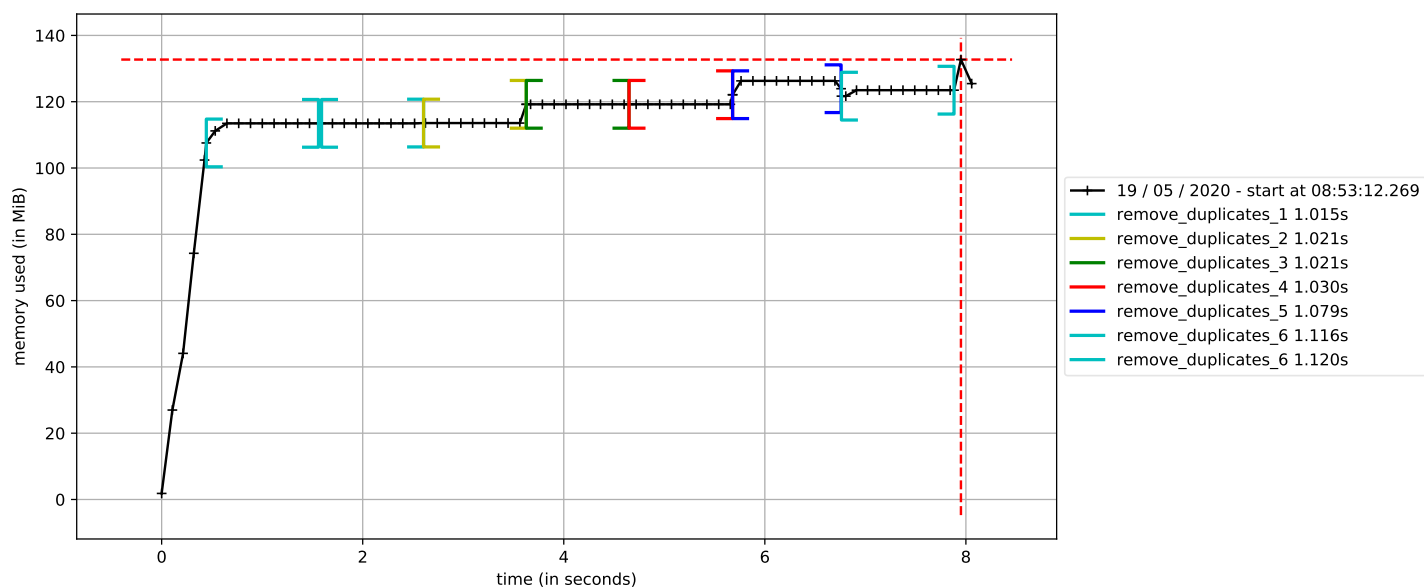
El código C está correctamente comentado, por lo que no veo necesario divagar mucho en el funcionamiento del mismo, pues es un tema que está mas relacionado con la algoritmia, y no con la esencia de esta práctica.

6. Resultados del profiling

6.1. Memory profiling

A continuación, en la siguiente página, muestro la gráfica resultante del análisis del uso de memoria realizado por las funciones python de juguete y el wrapper. En ella se pueden contrastar los resultados junto a la explicación que he realizado en el apartado anterior. Este gráfico ha sido generado a través del archivo `/doc/profiling/mpprofile_20200519085312.dat`

/home/victorio/miniconda3/bin/python main.py ../../test/numbers.txt results.txt results.pdf



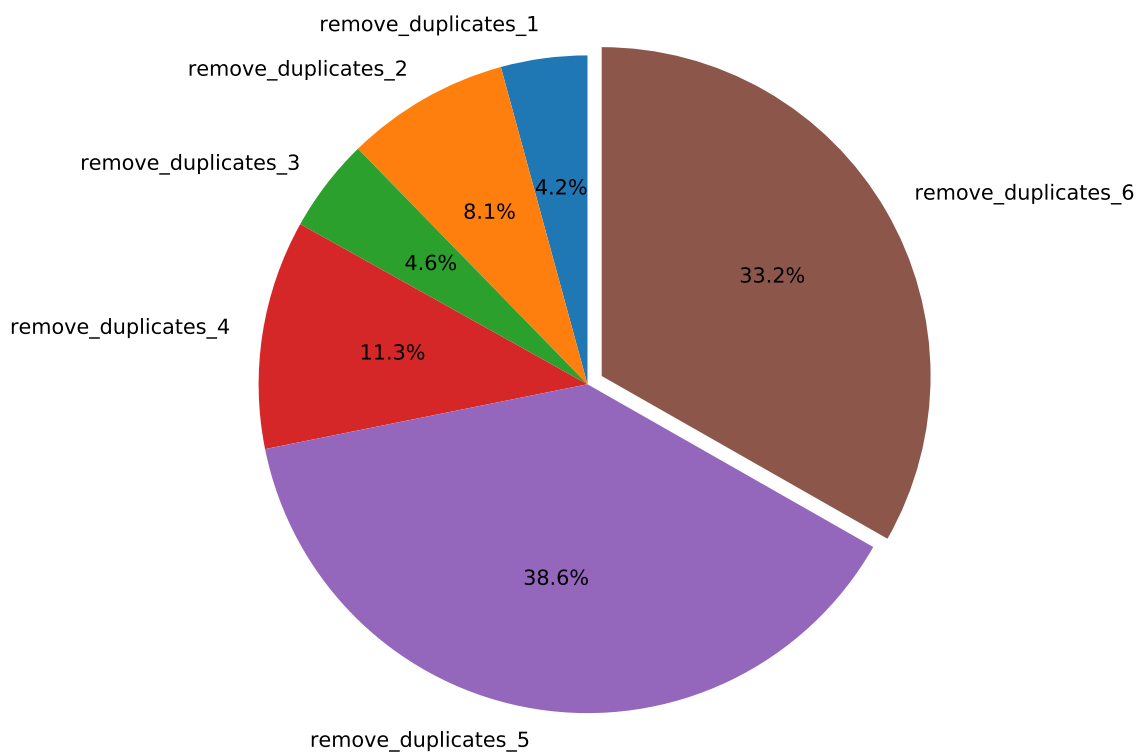
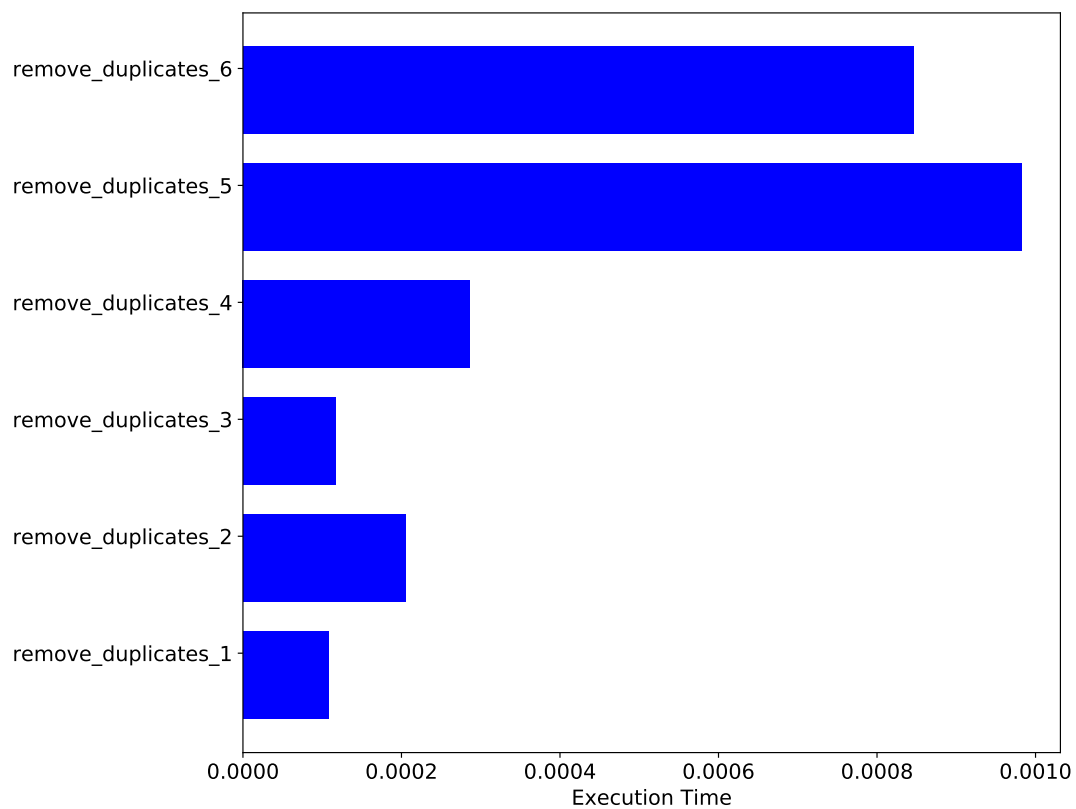
6.1.1. Linear profiling y gráficas generadas con Matplotlib

El profiling del tiempo de ejecución lo he realizado con el módulo *linear_profiling*, tal y como he comentado anteriormente. En el directorio **/doc/profiling** se puede encontrar su resultado, concretamente en el fichero **main.py.lprof**. Para terminar con el contraste de resultados, voy a incluir una serie de figuras (creadas con Matplotlib) compuestas por dos gráficos, uno de barras y otro de tarta, como ya comenté anteriormente.

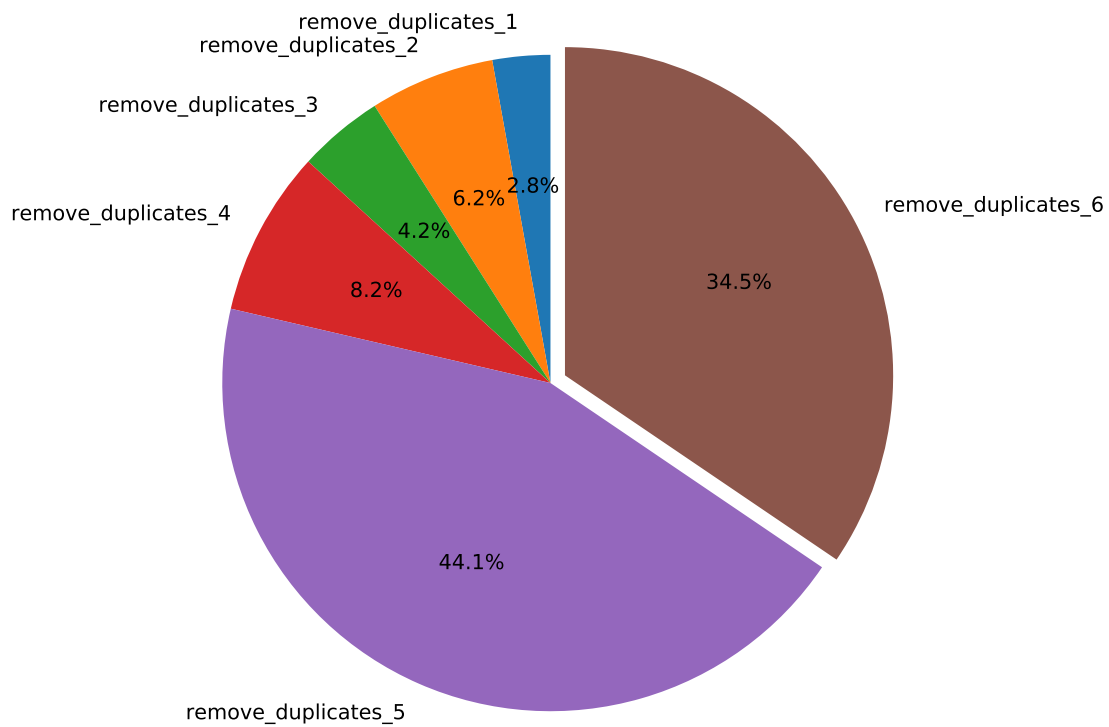
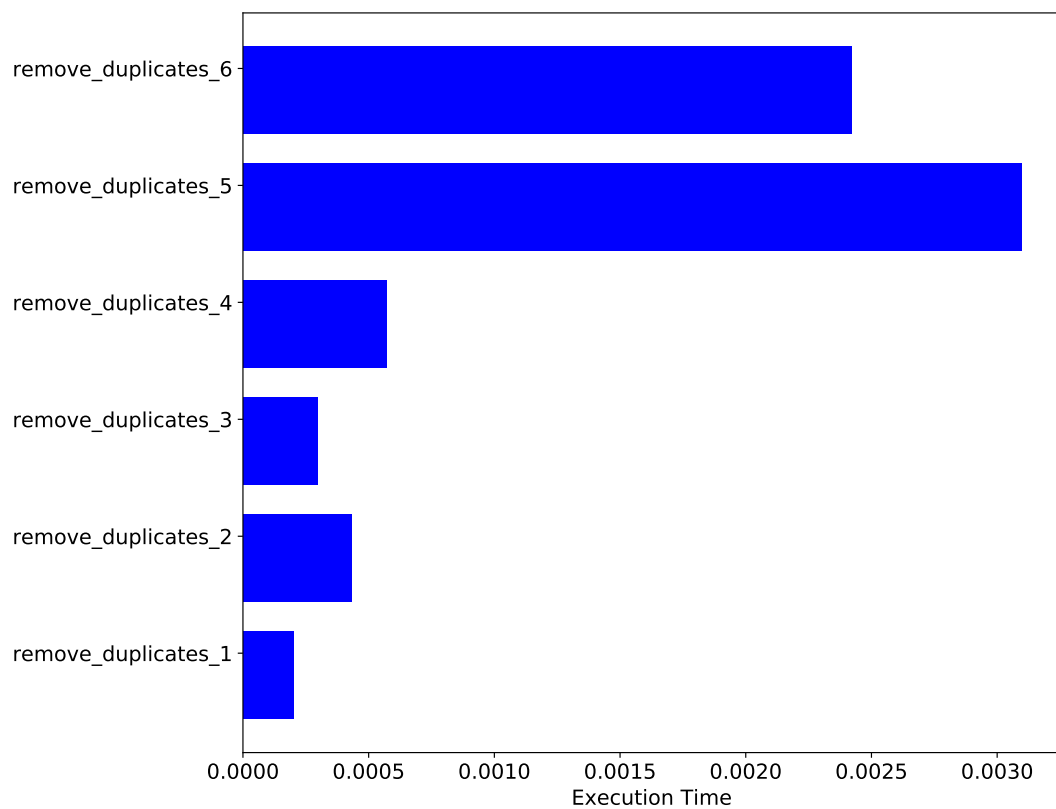
Tal y como se comenta en el enunciado, el análisis lo he realizado, además de usando el módulo *linear_profiling*, generando gráficas que muestran el tiempo usado por cada técnica para los distintos tamaños de datos, de 2000 en 2000, hasta llegar al tamaño total de la lista de elementos del fichero de entrada.

Pd: el siguiente contenido lo almacena el programa en el fichero correspondiente al tercer argumento que le pasa el usuario.

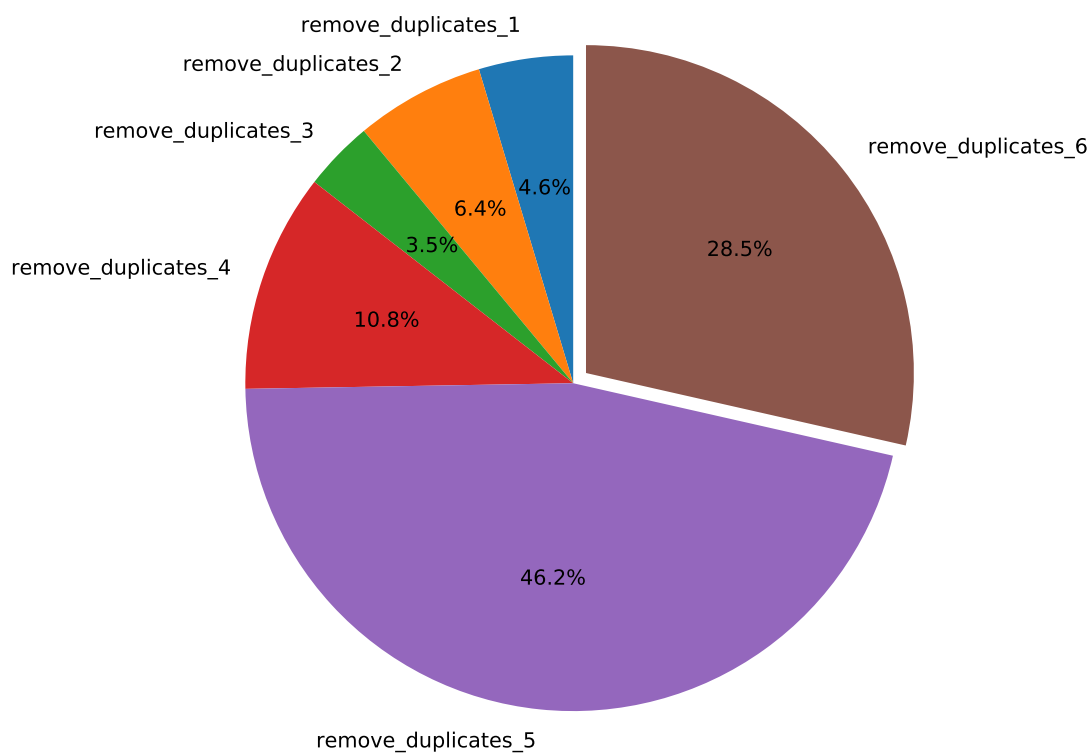
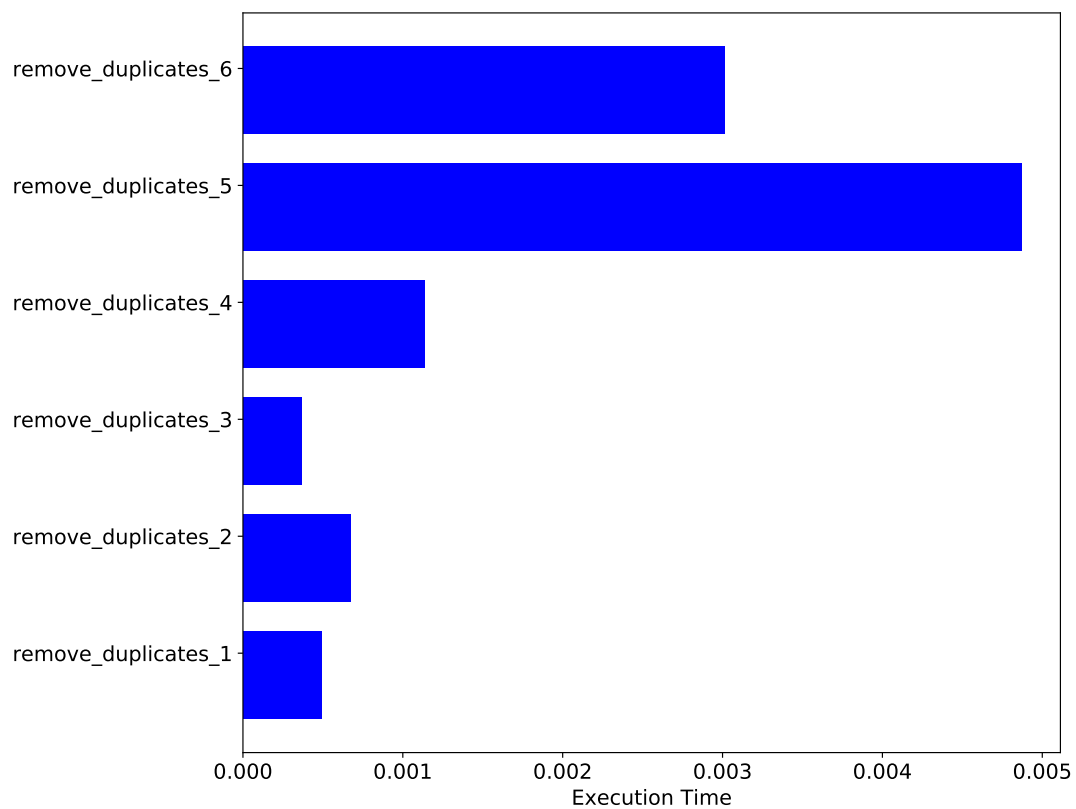
TIMES FOR THE FIRST 2000 NUMBERS



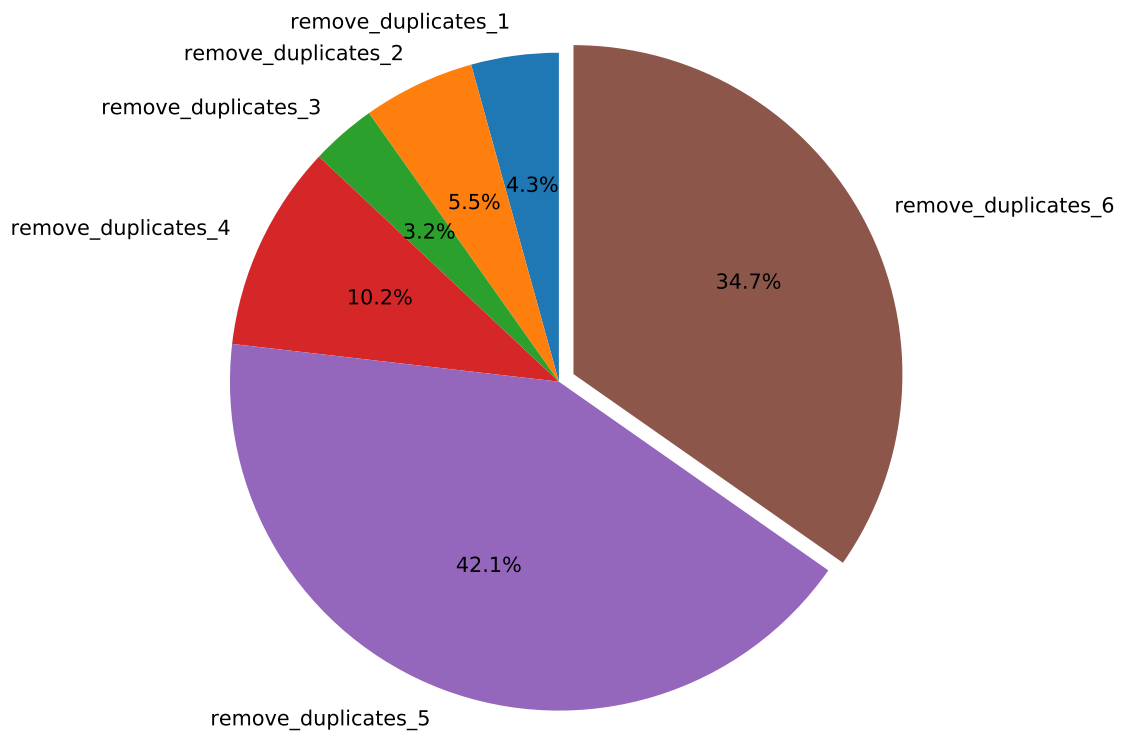
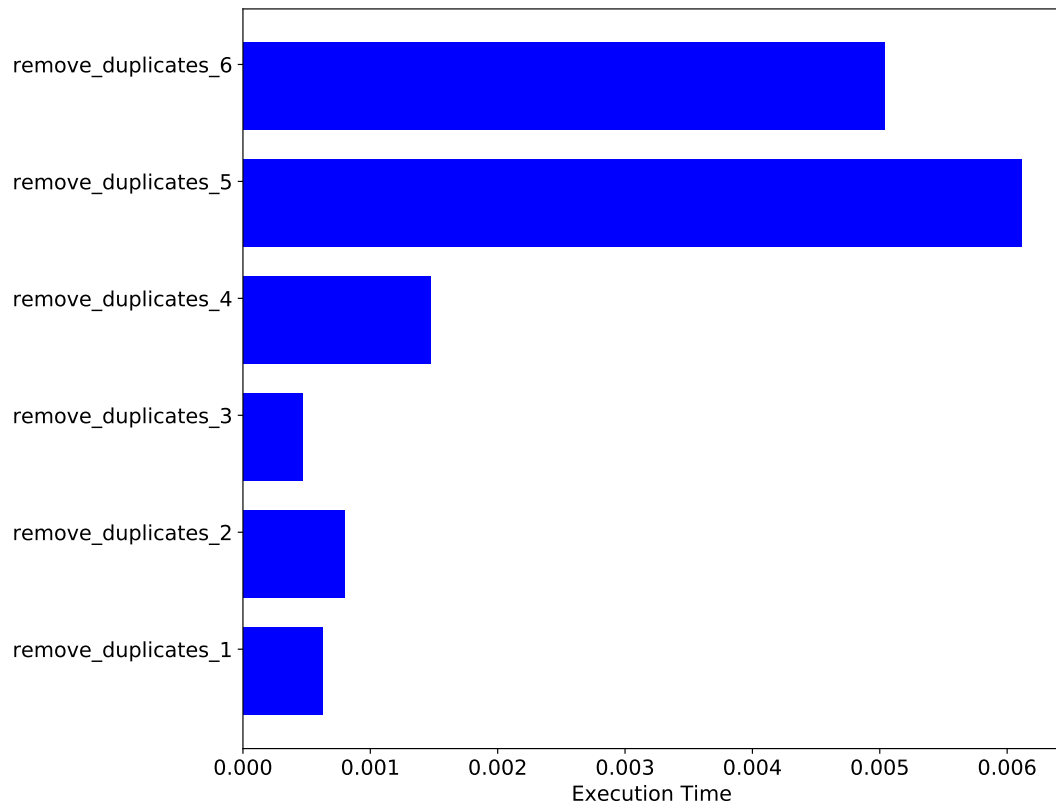
TIMES FOR THE FIRST 4000 NUMBERS



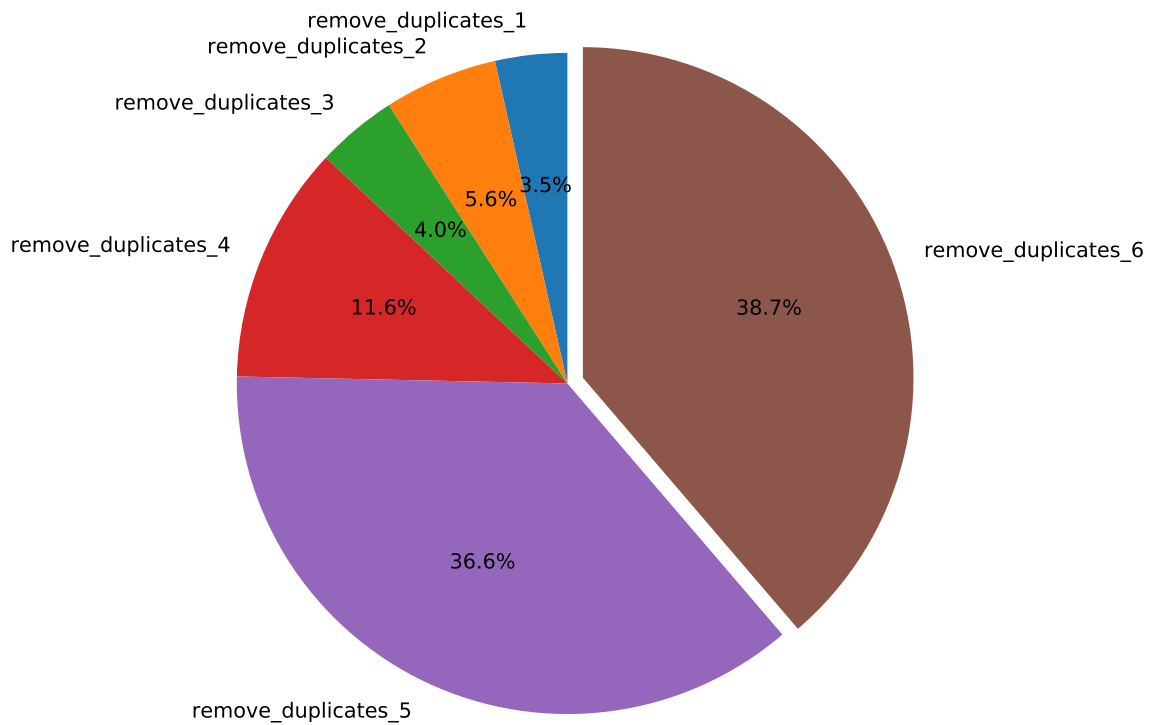
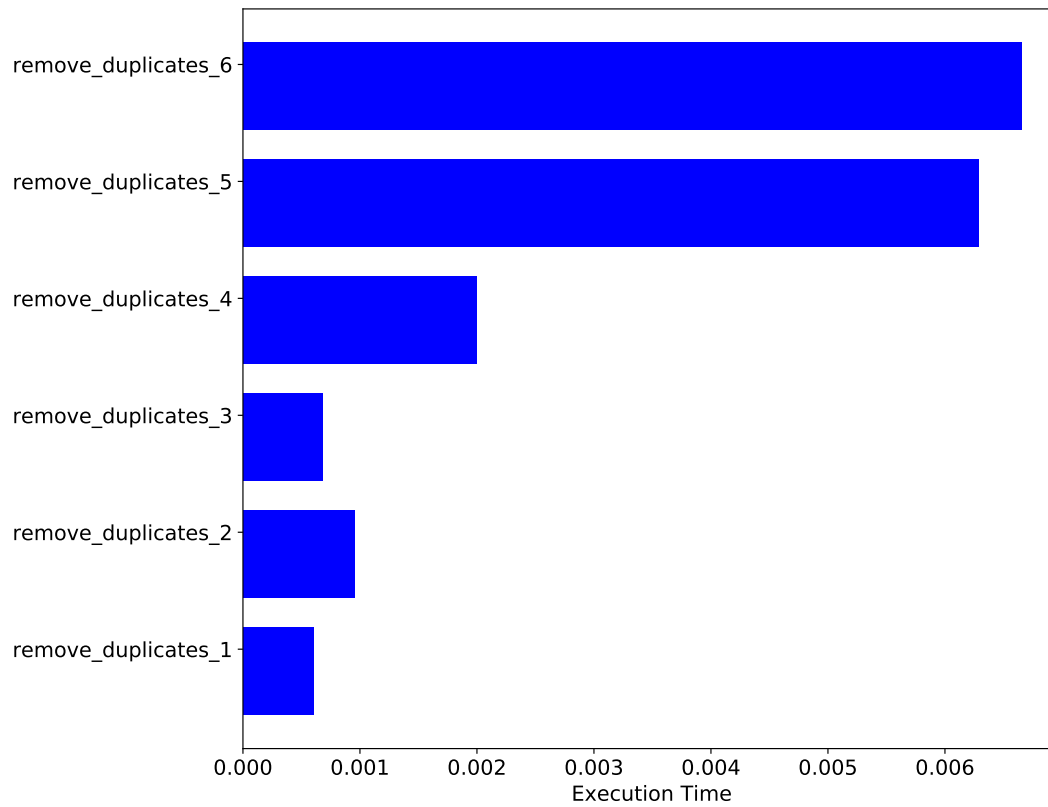
TIMES FOR THE FIRST 6000 NUMBERS



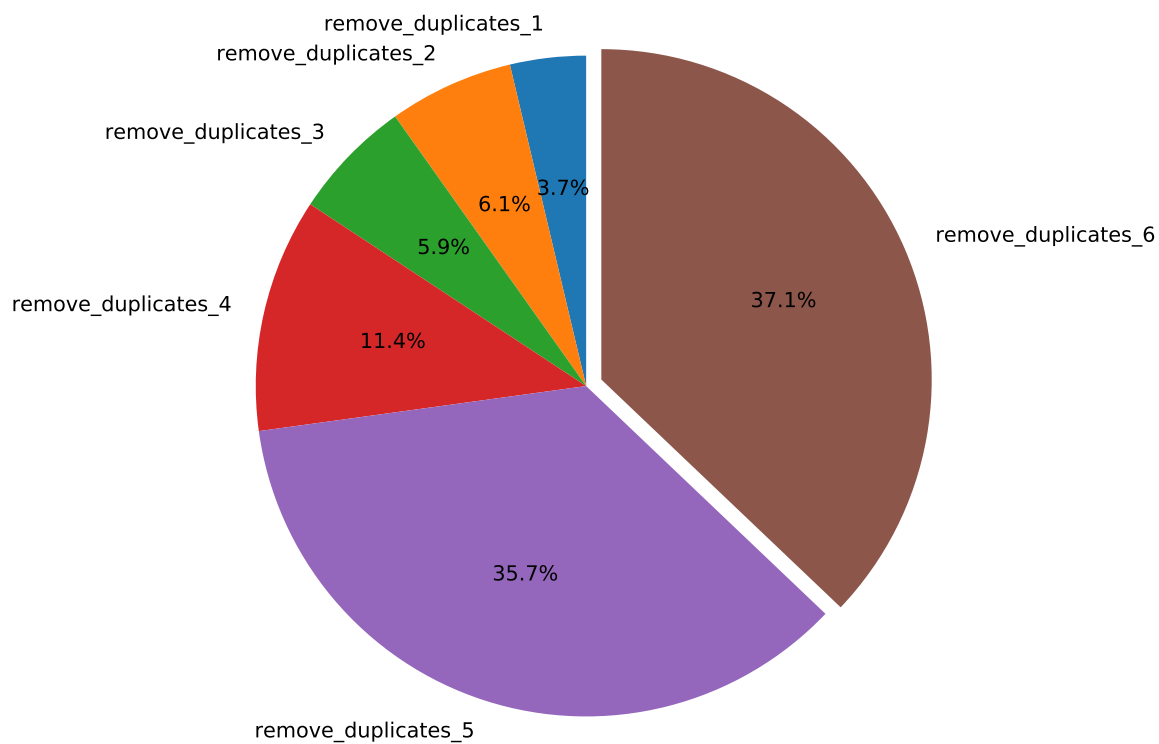
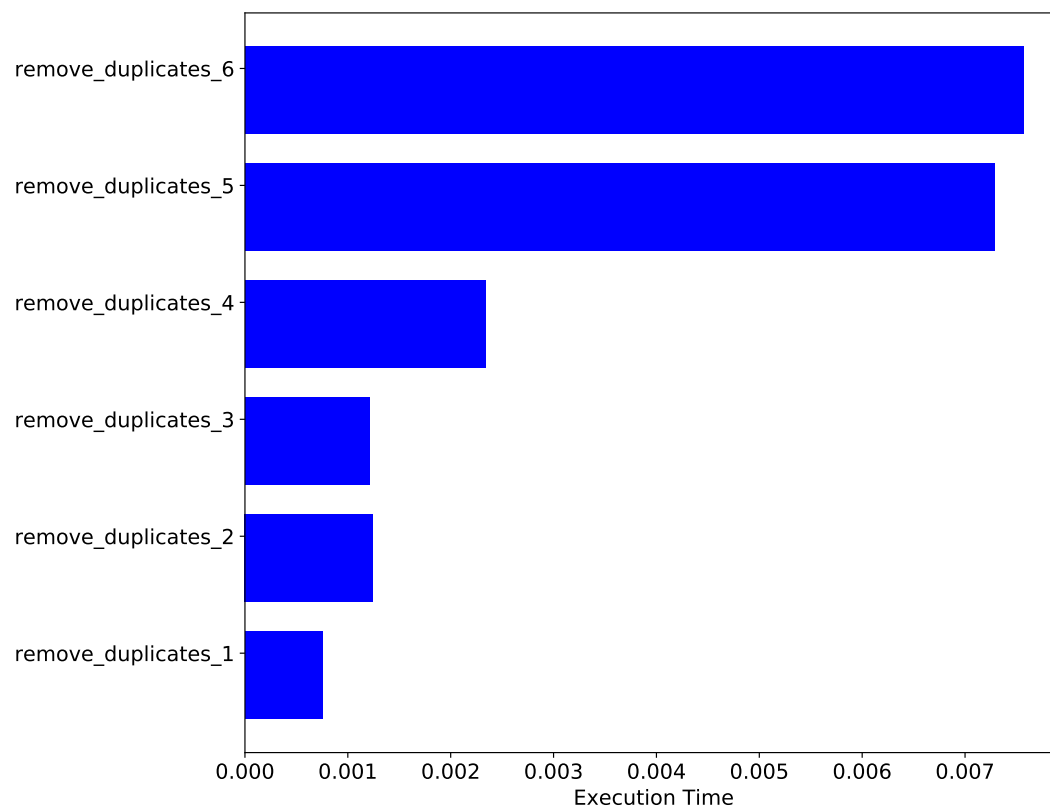
TIMES FOR THE FIRST 8000 NUMBERS



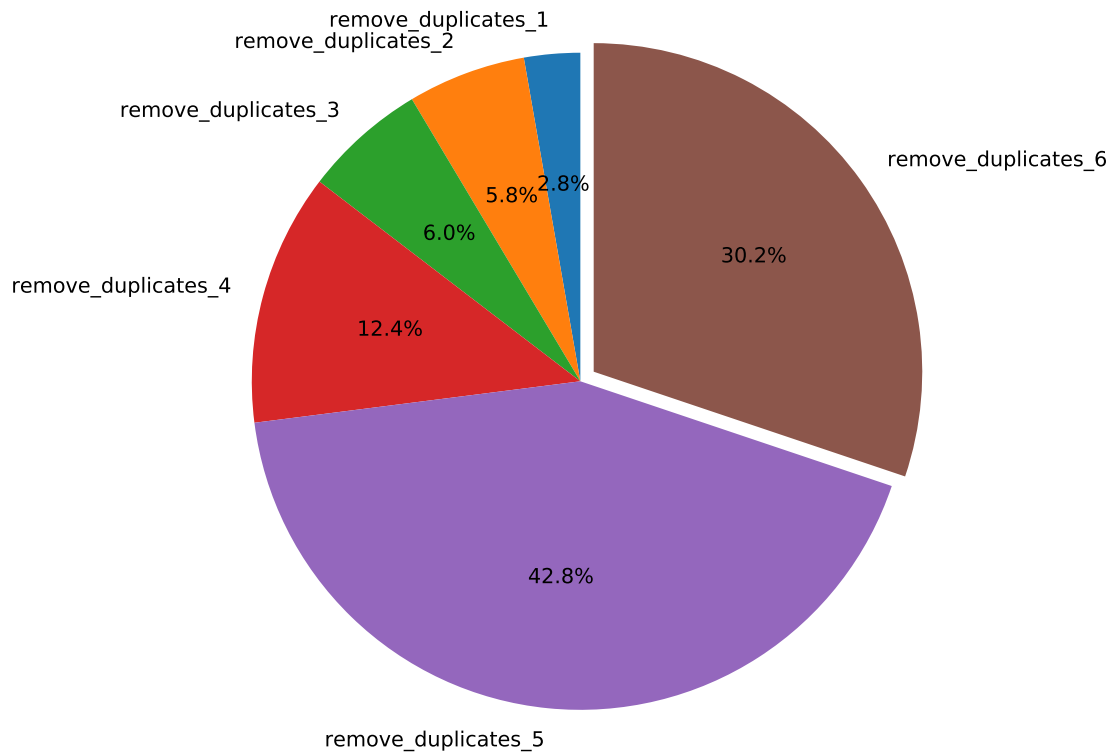
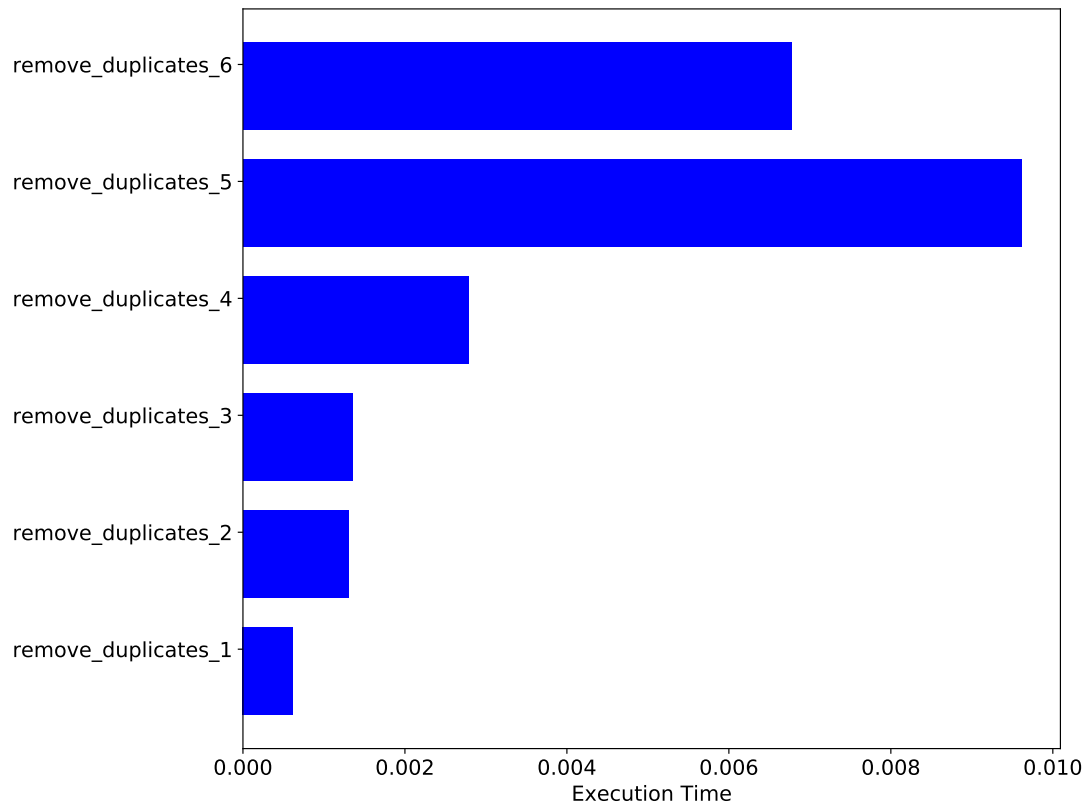
TIMES FOR THE FIRST 10000 NUMBERS



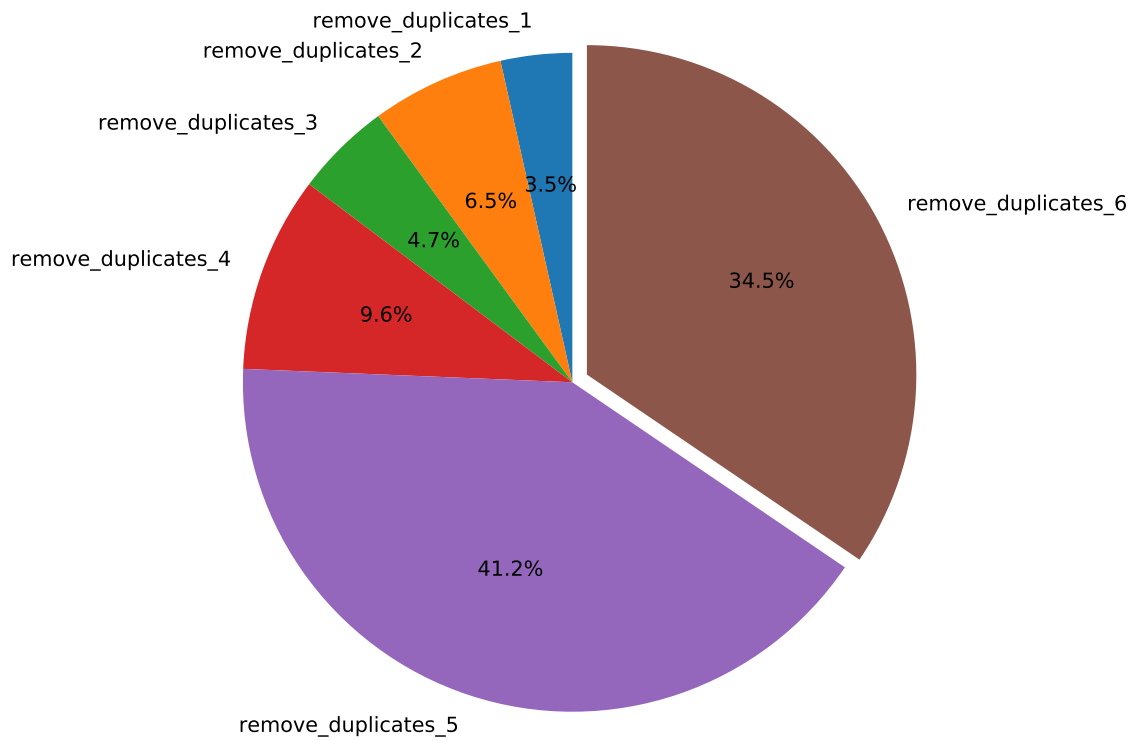
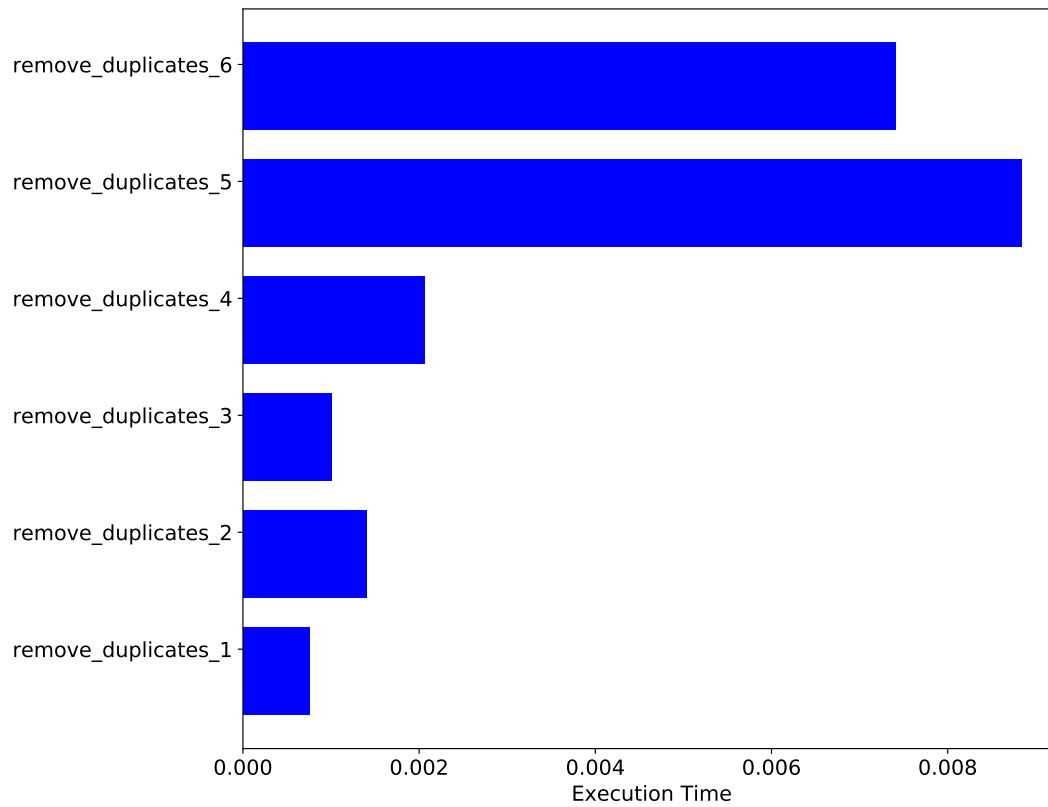
TIMES FOR THE FIRST 12000 NUMBERS



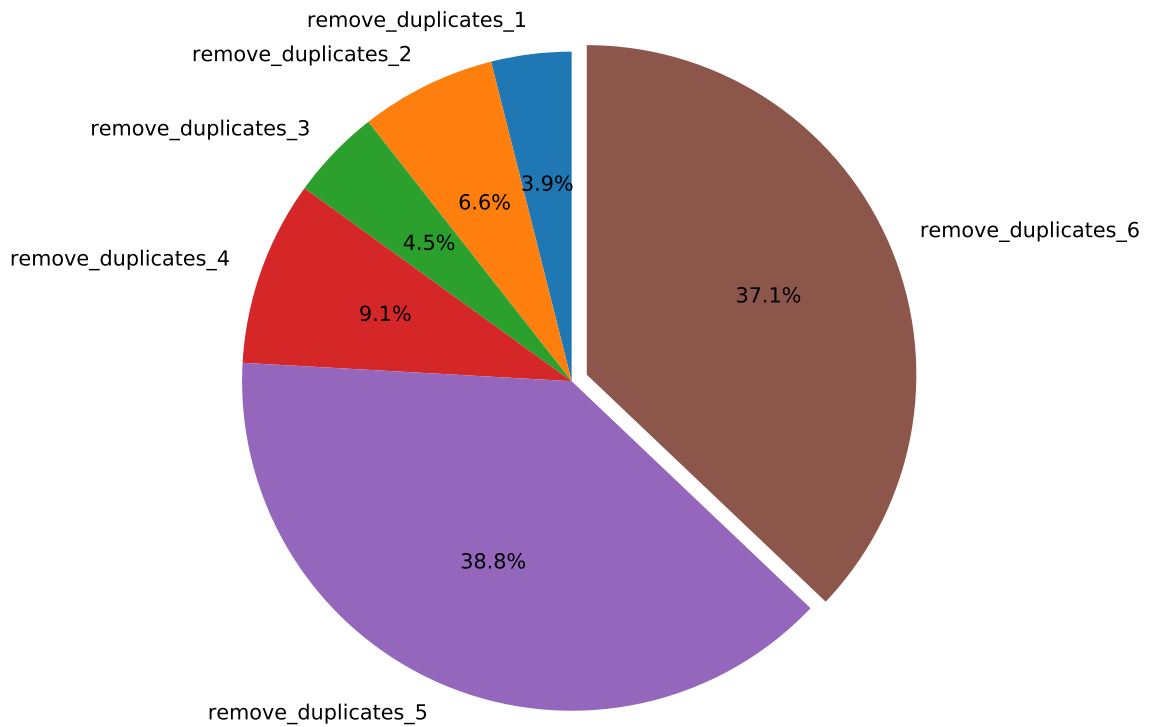
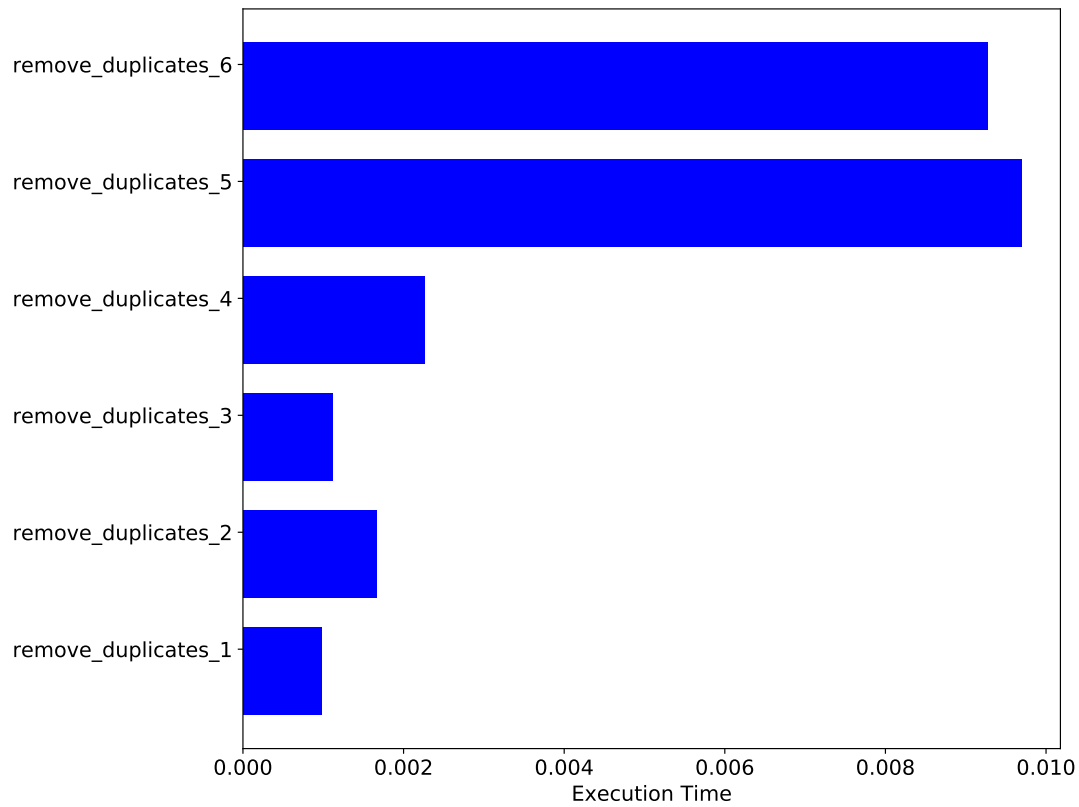
TIMES FOR THE FIRST 14000 NUMBERS



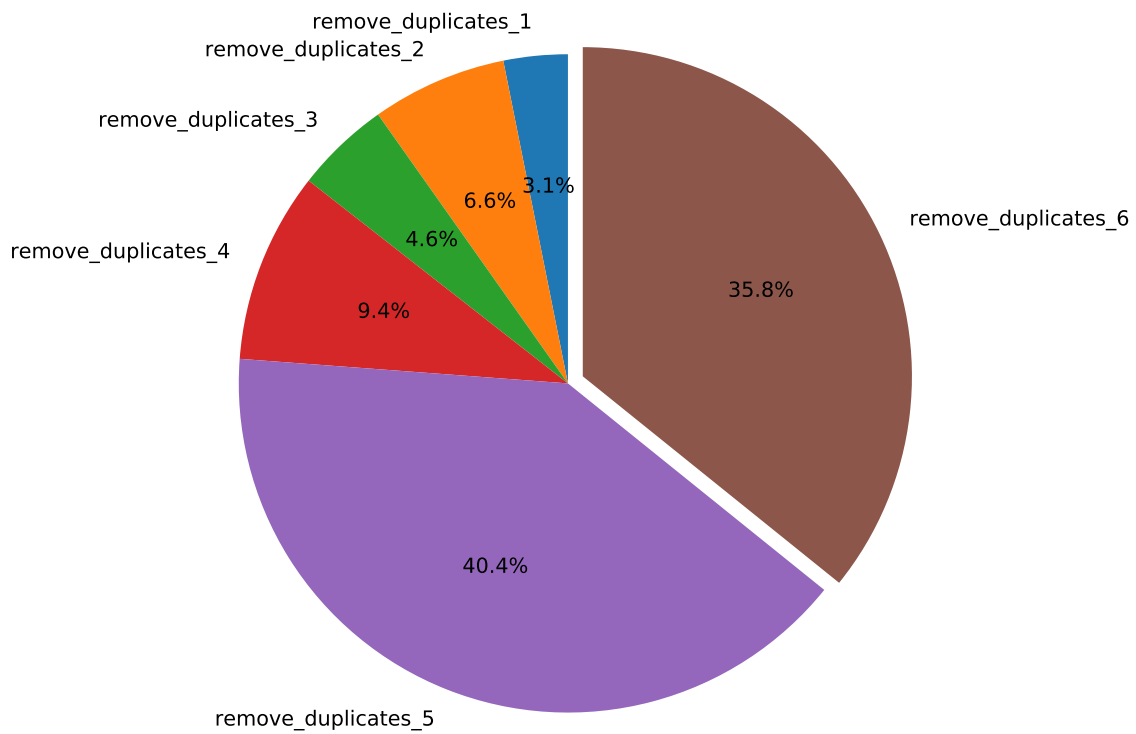
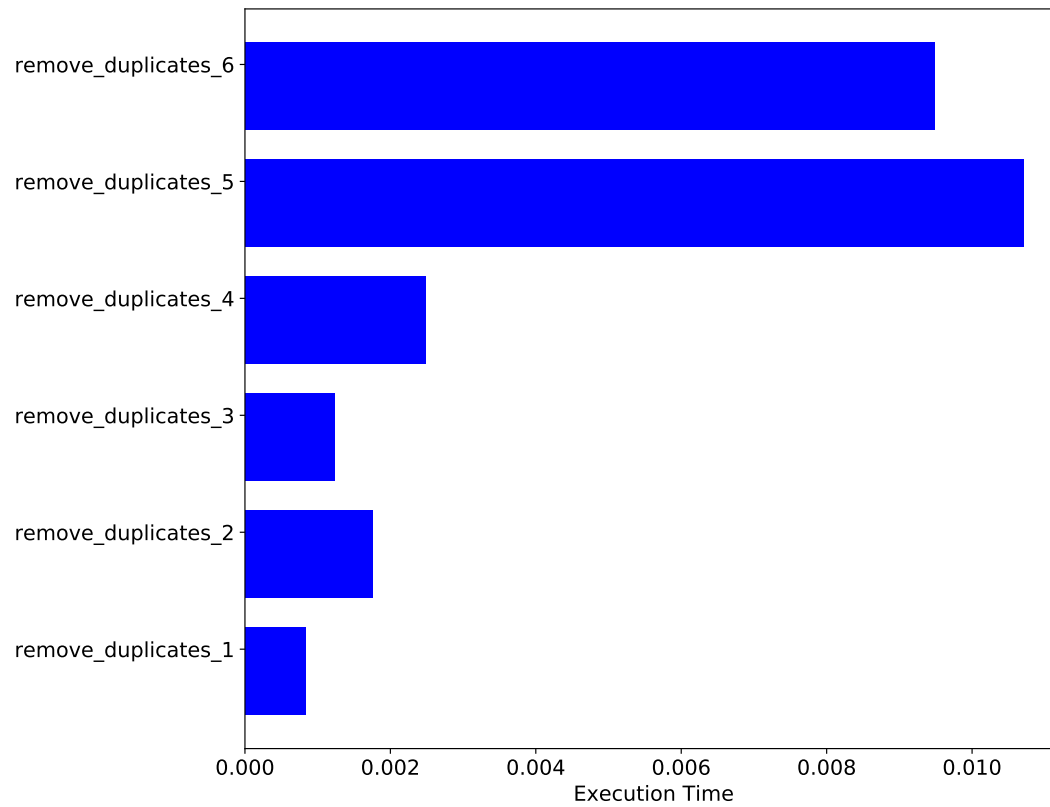
TIMES FOR THE FIRST 16000 NUMBERS



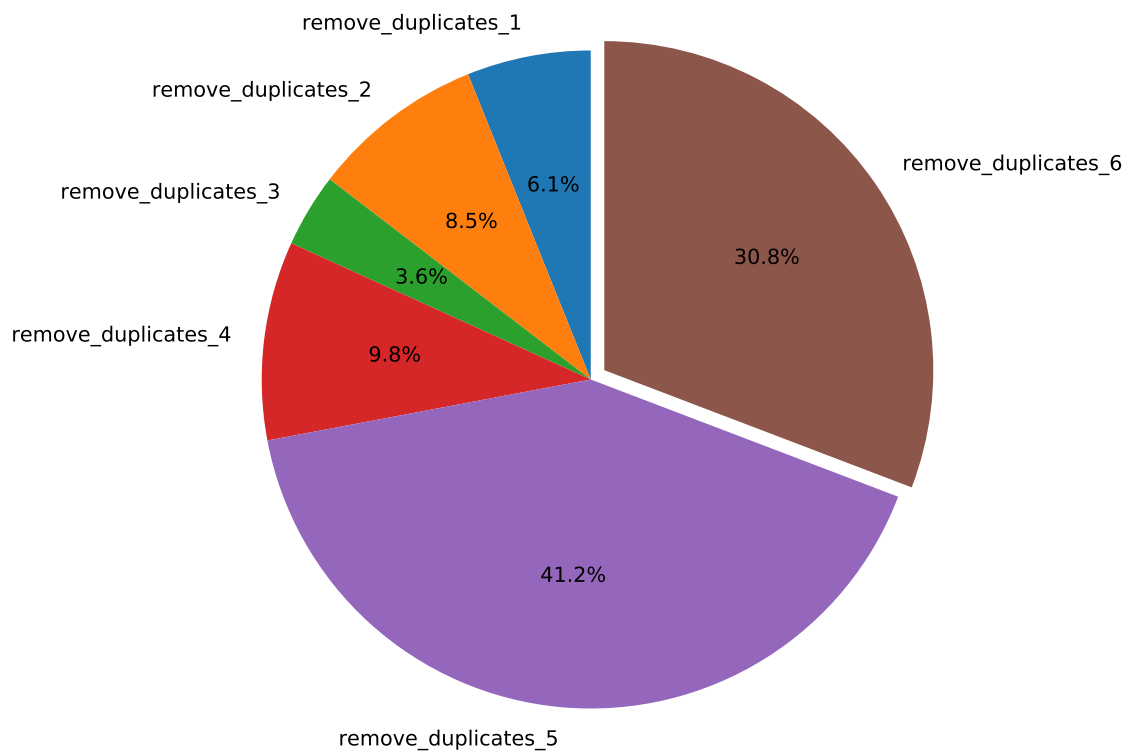
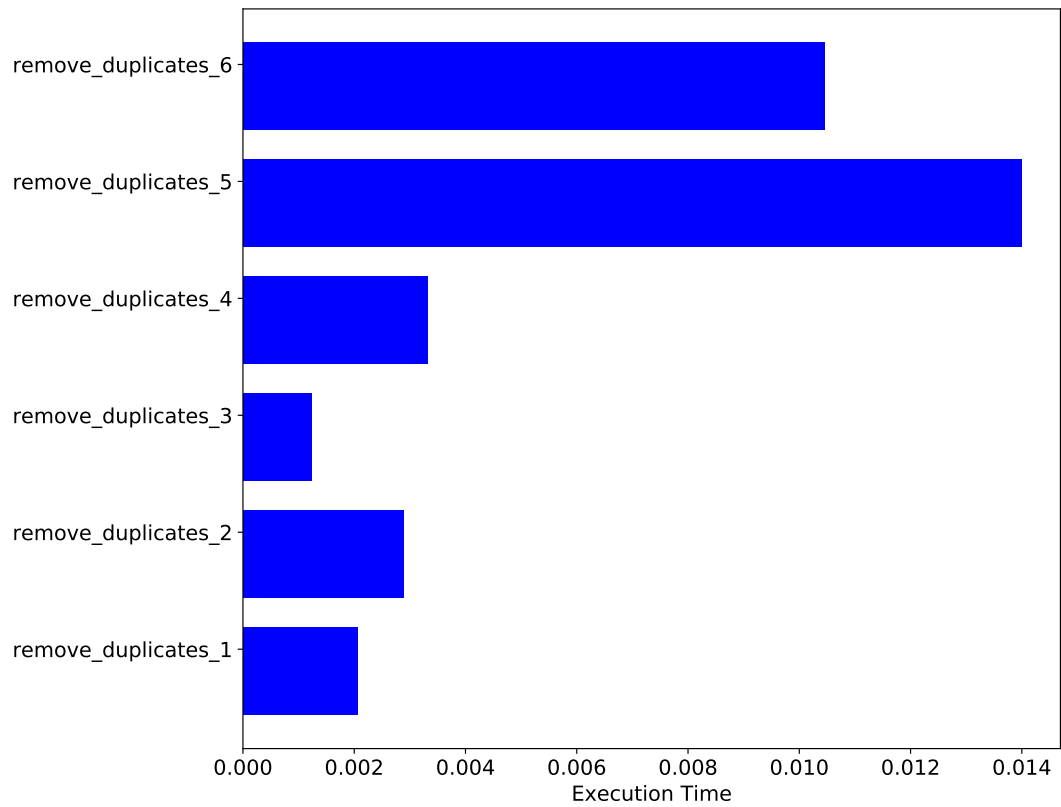
TIMES FOR THE FIRST 18000 NUMBERS



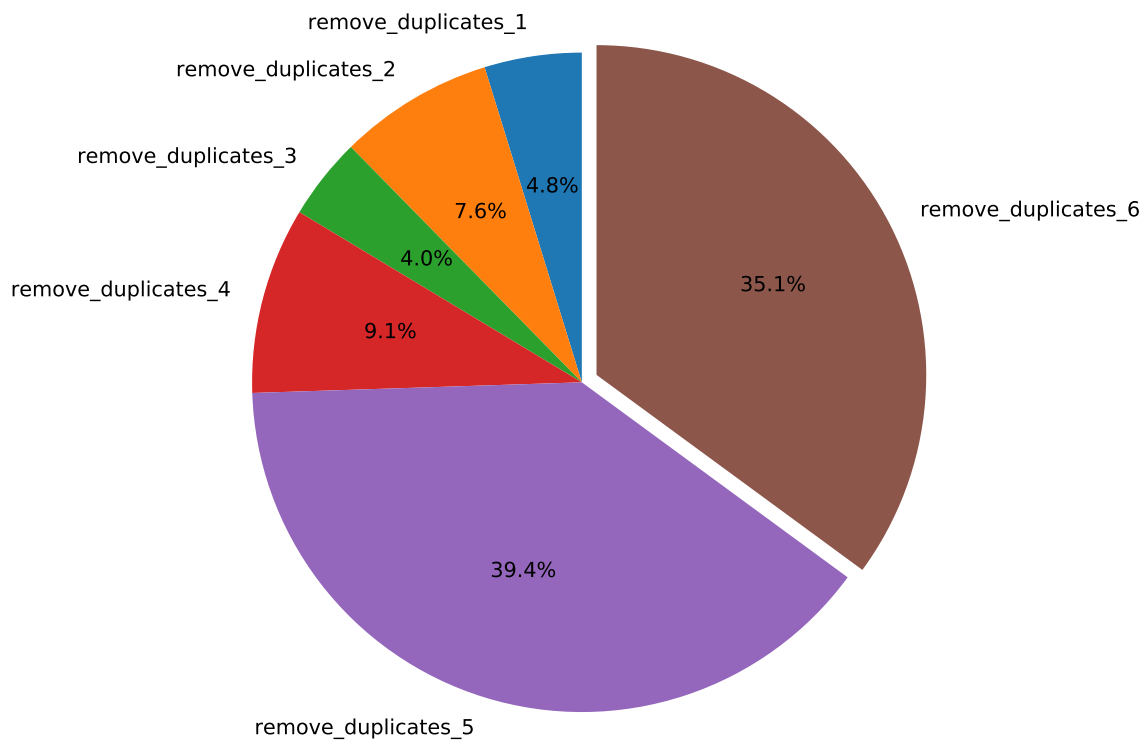
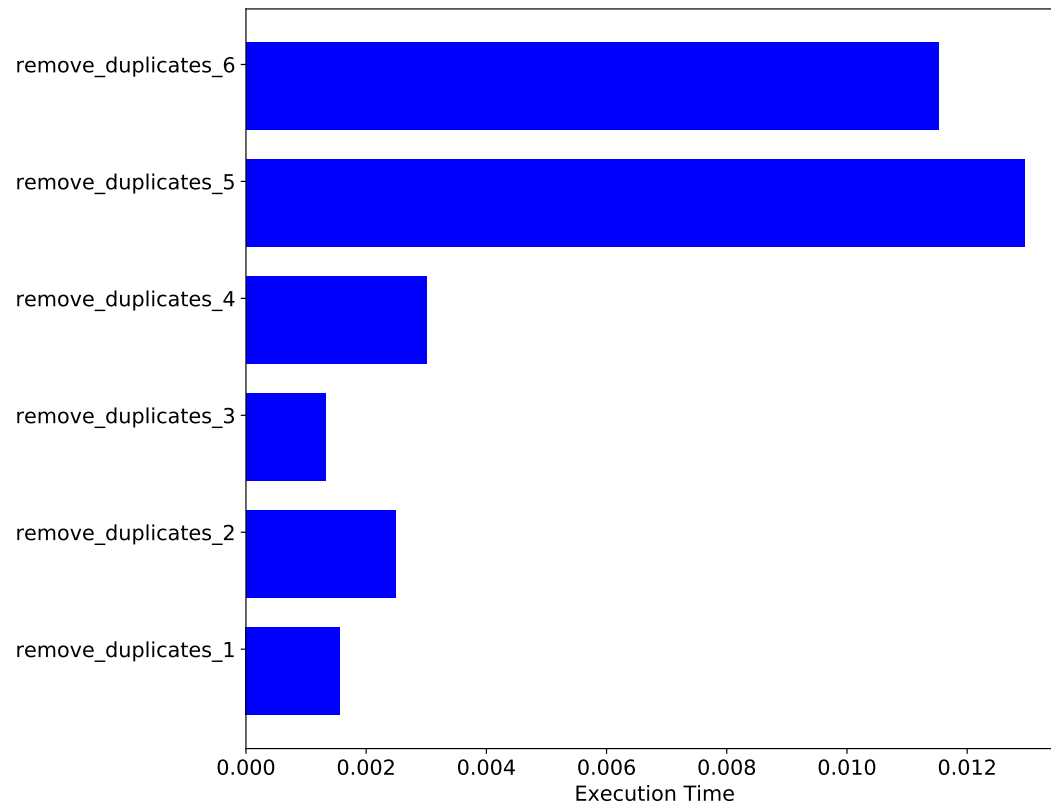
TIMES FOR THE FIRST 20000 NUMBERS



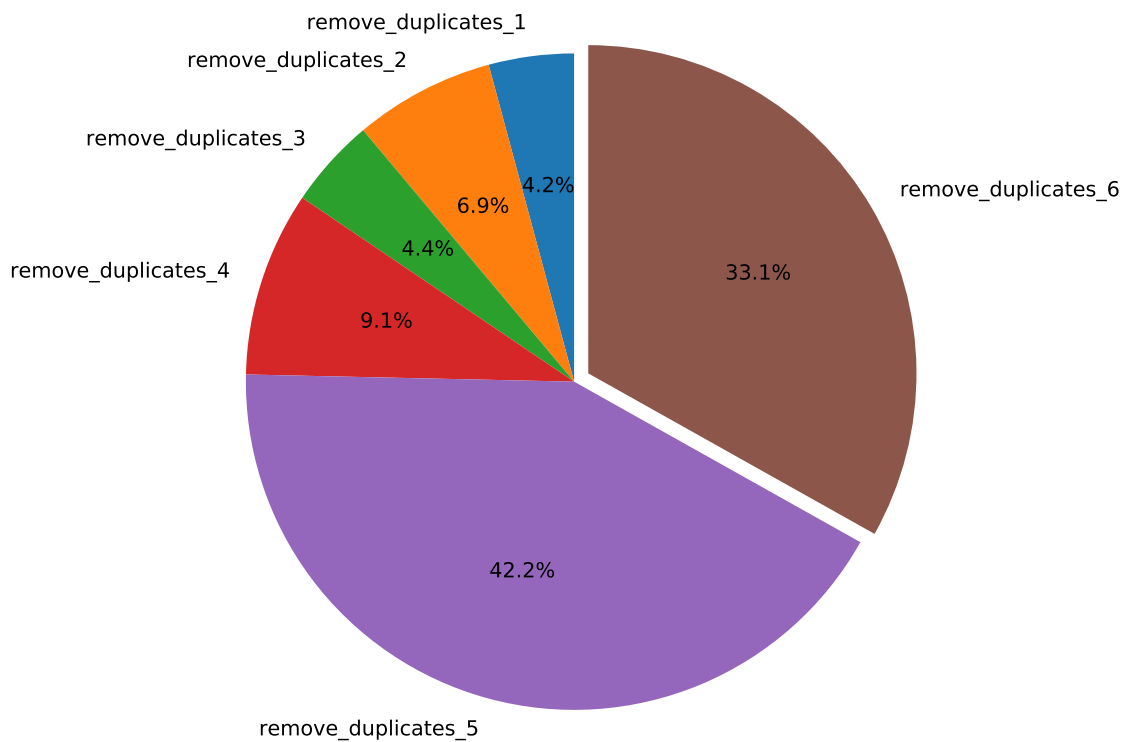
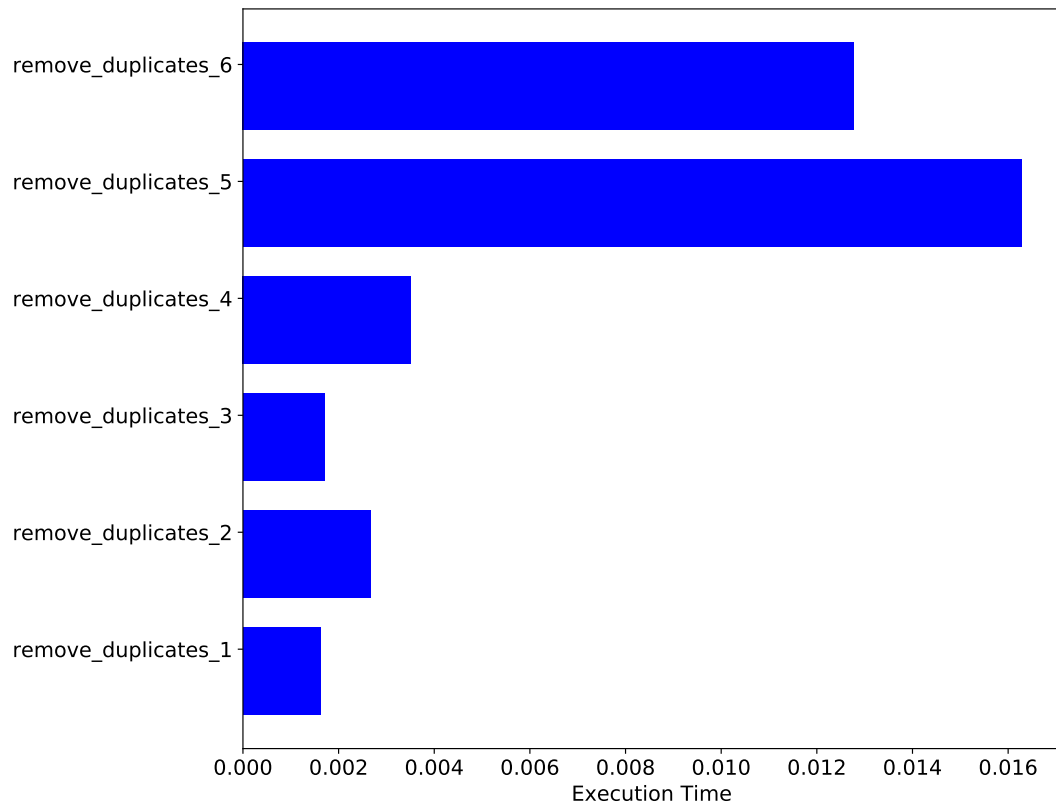
TIMES FOR THE FIRST 22000 NUMBERS



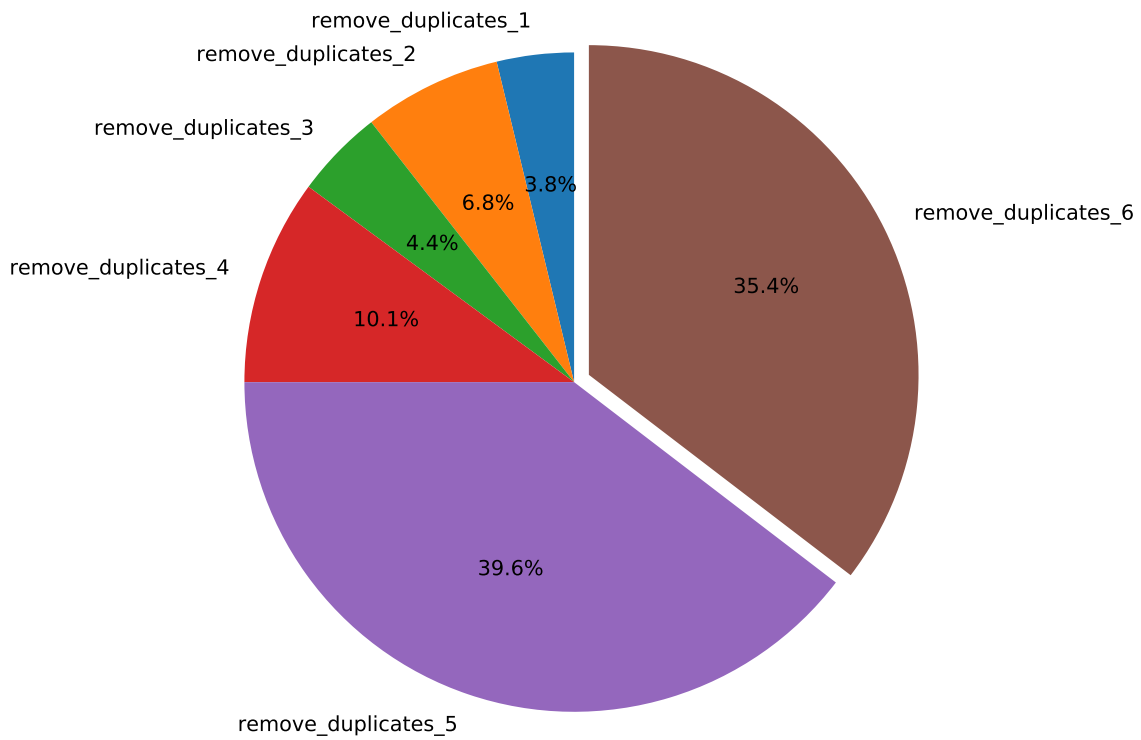
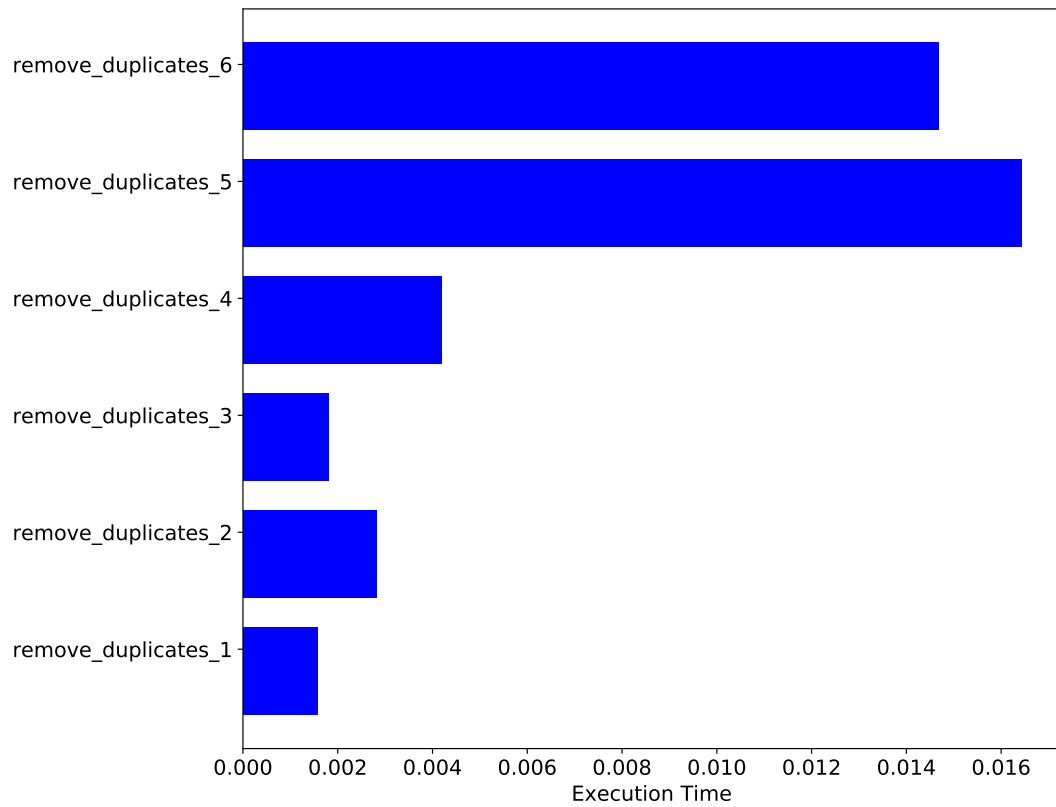
TIMES FOR THE FIRST 24000 NUMBERS



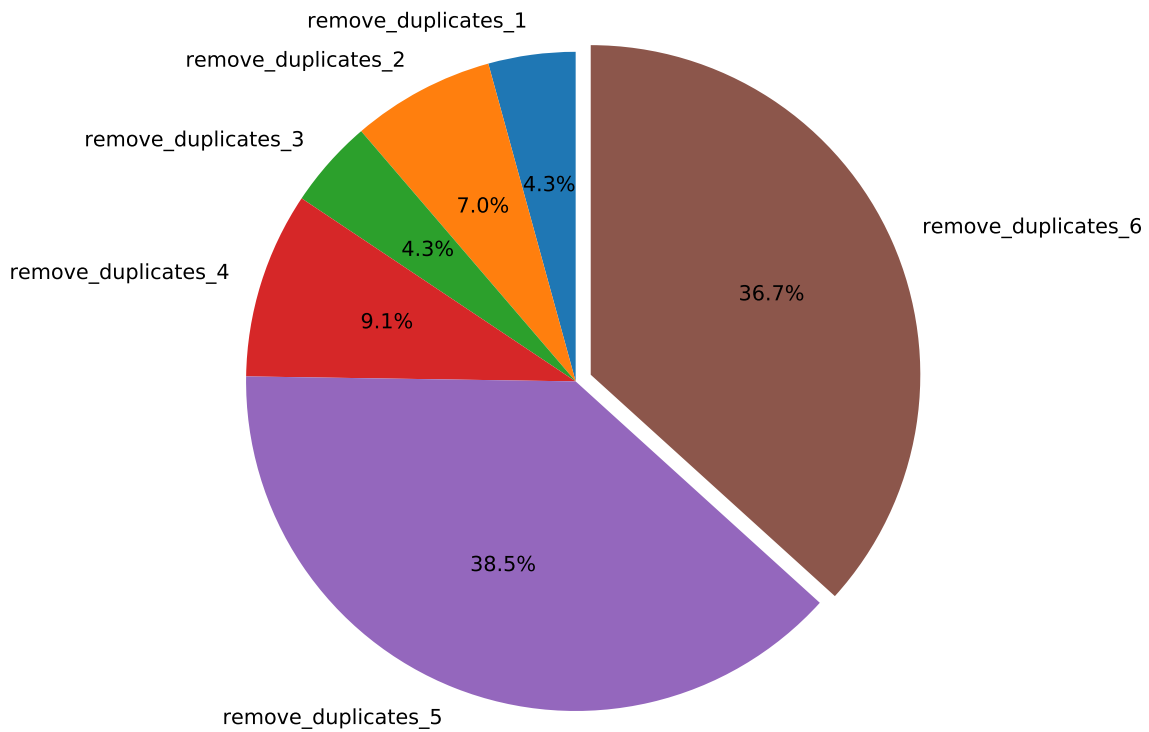
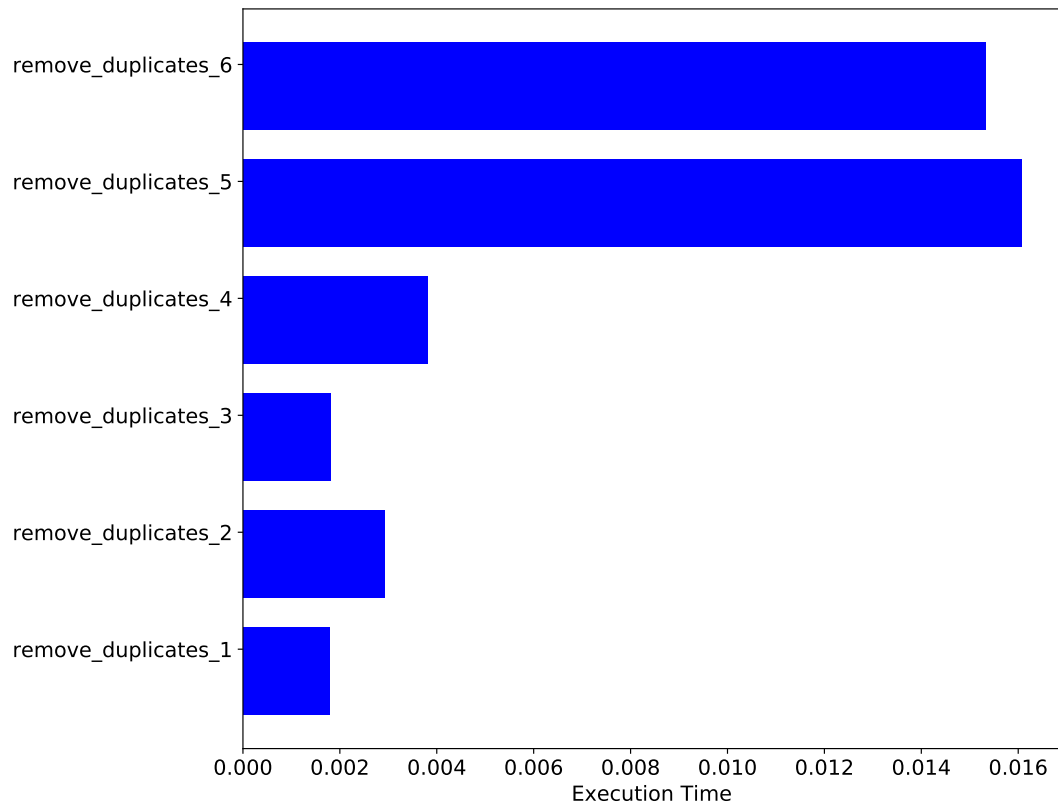
TIMES FOR THE FIRST 26000 NUMBERS



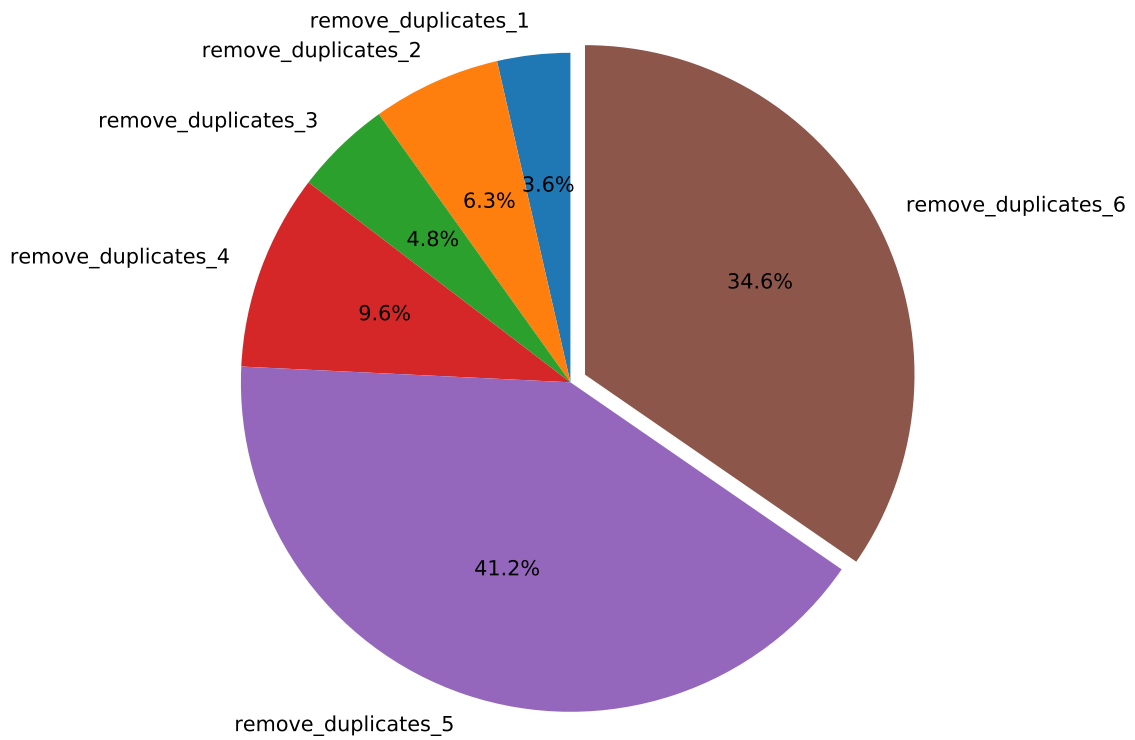
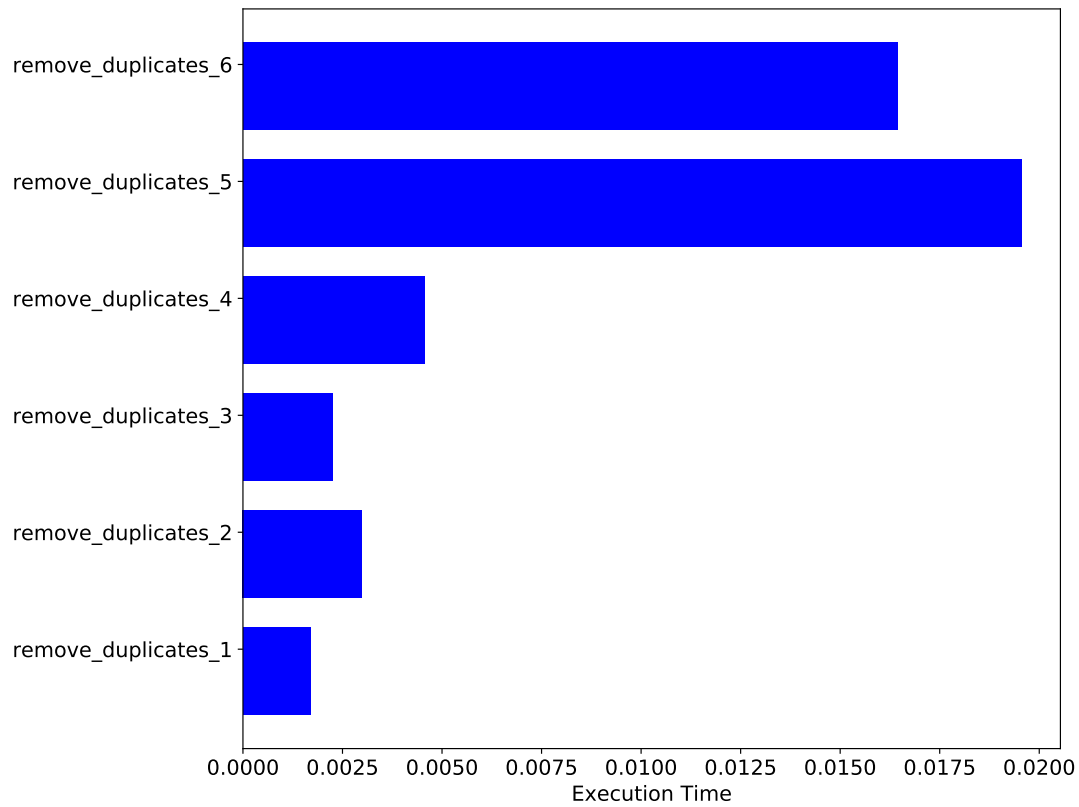
TIMES FOR THE FIRST 28000 NUMBERS



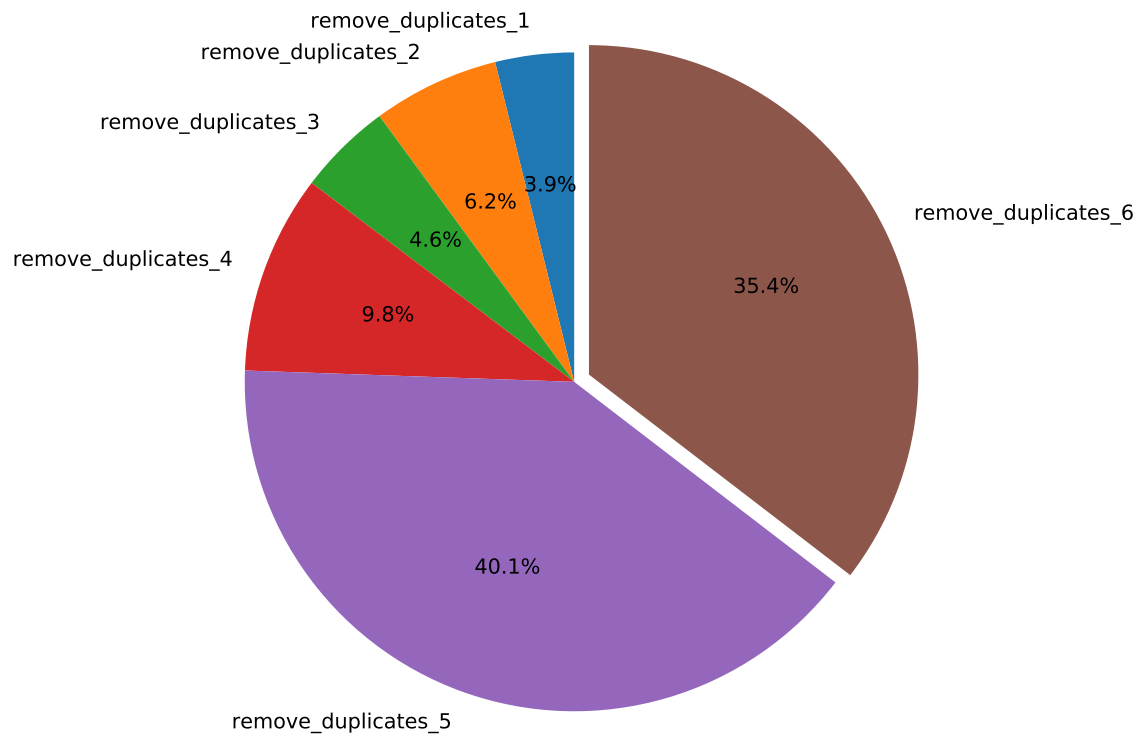
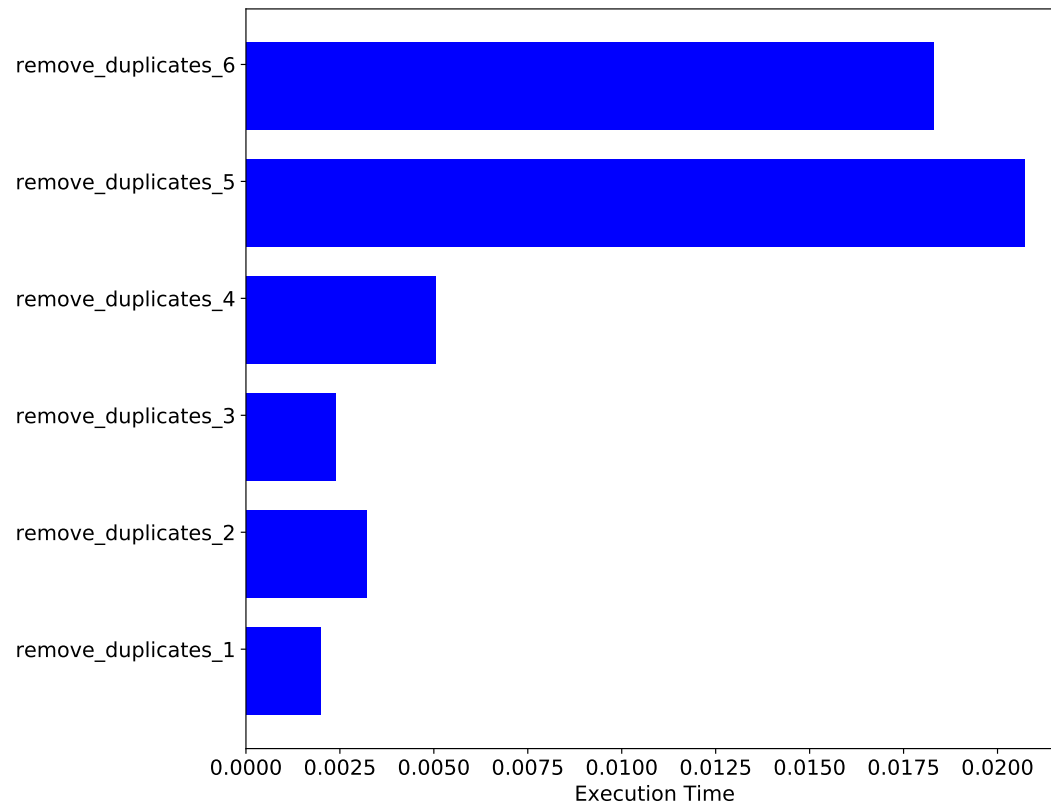
TIMES FOR THE FIRST 30000 NUMBERS



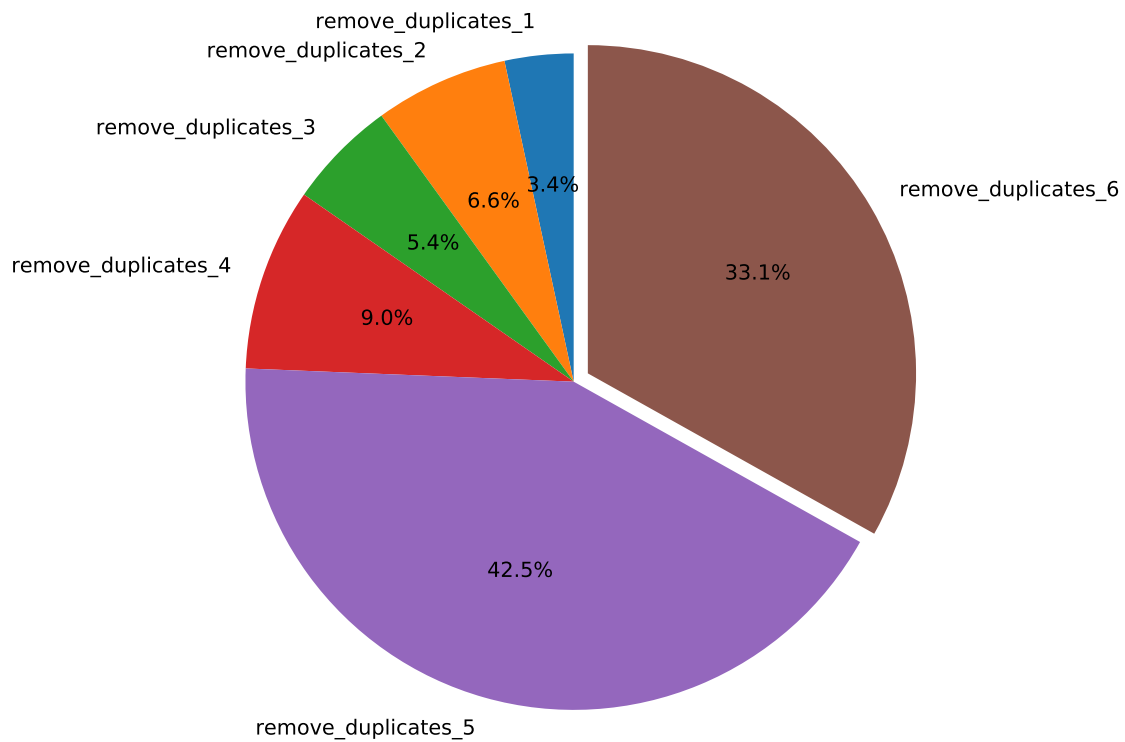
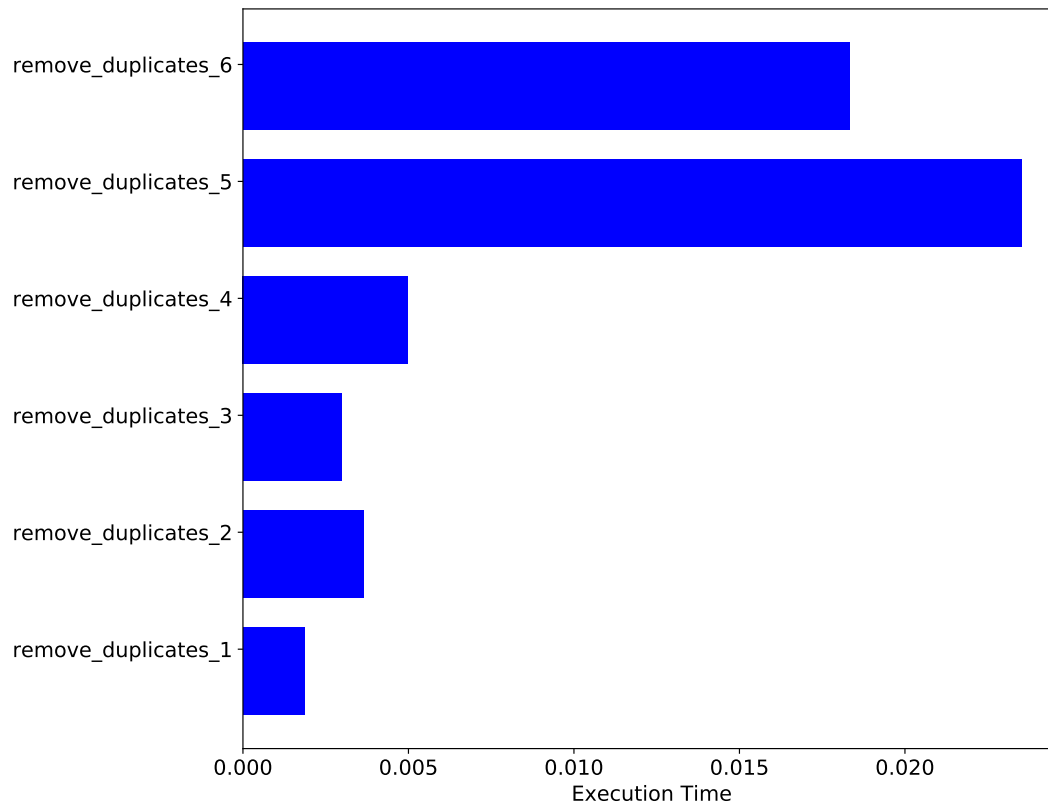
TIMES FOR THE FIRST 32000 NUMBERS



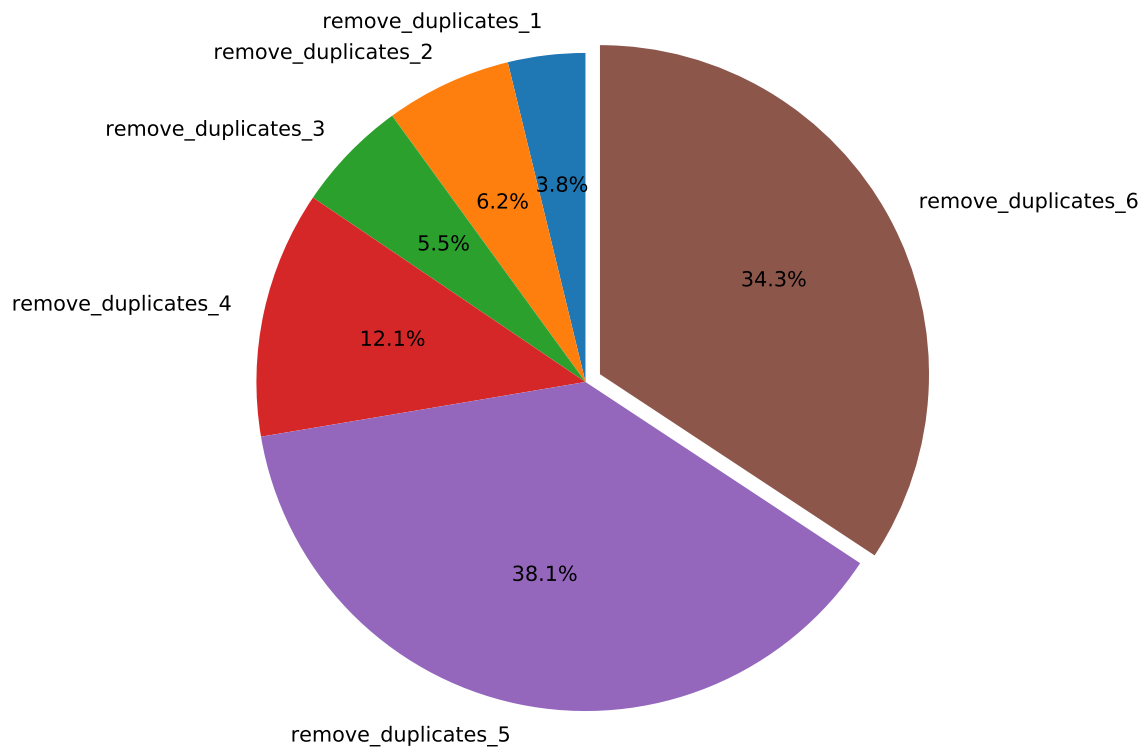
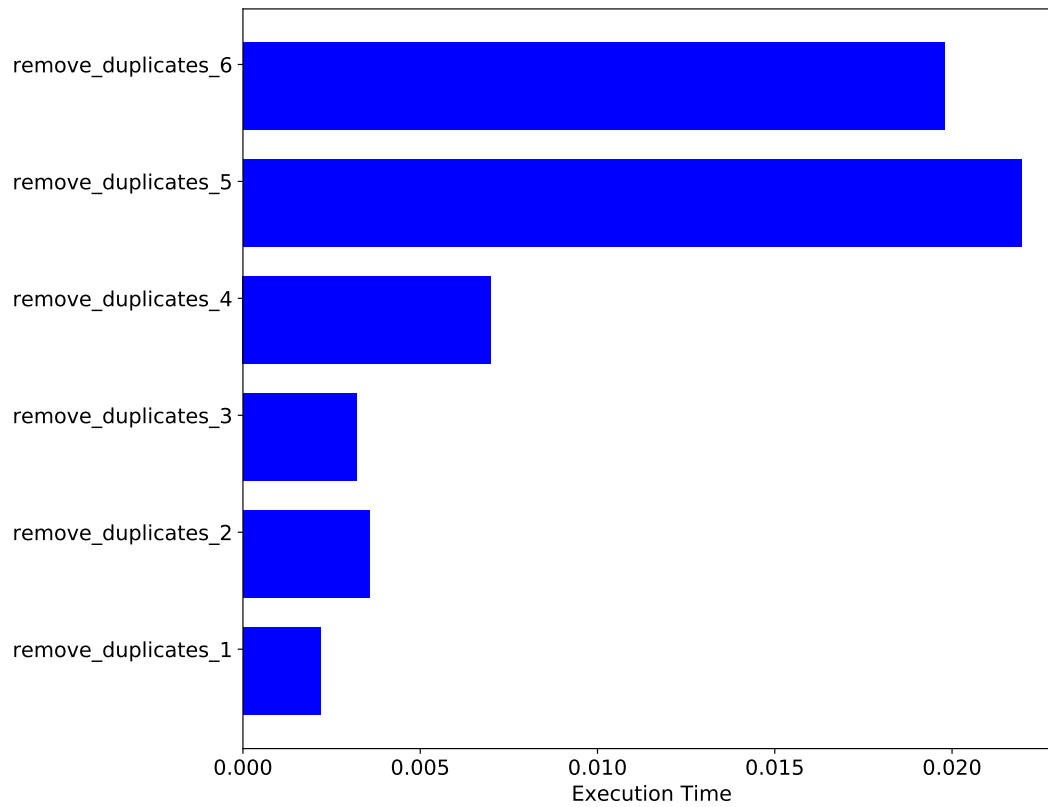
TIMES FOR THE FIRST 34000 NUMBERS



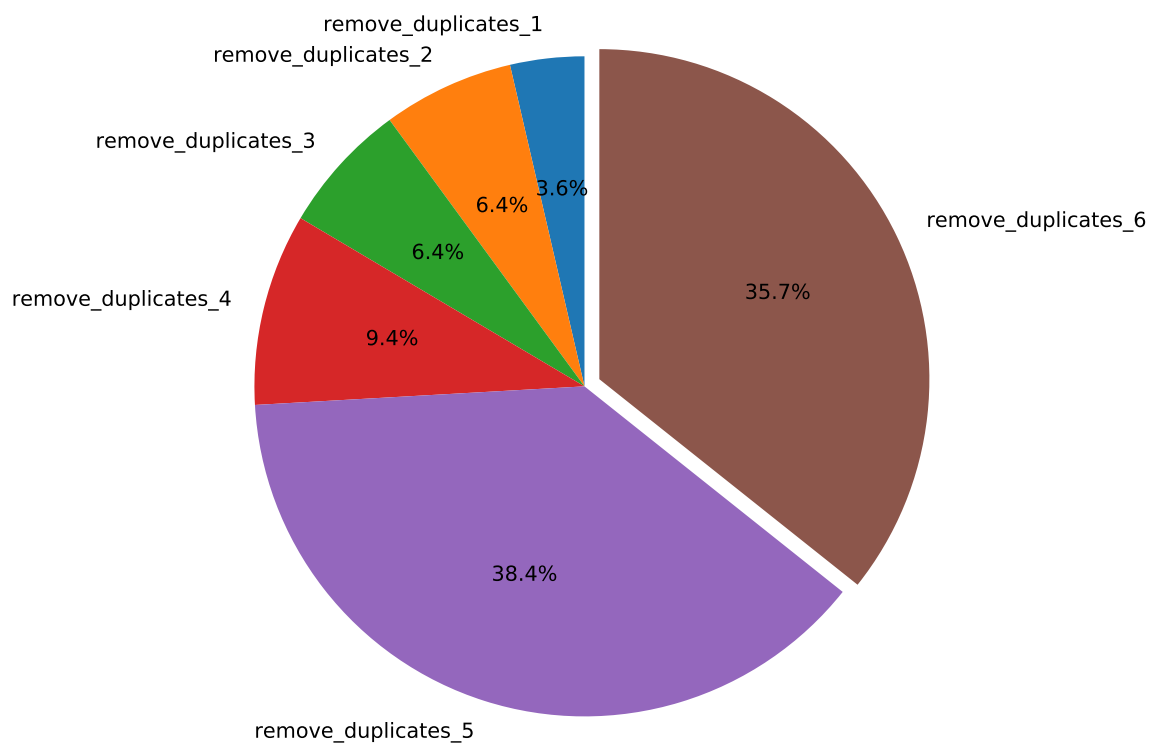
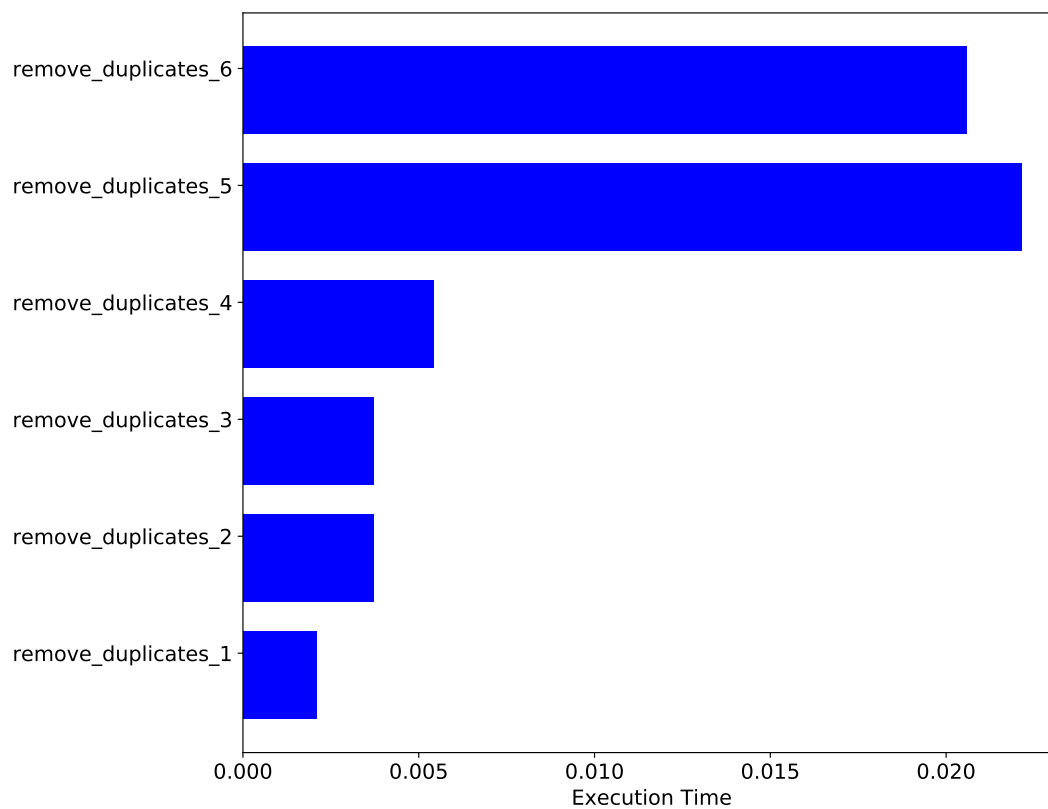
TIMES FOR THE FIRST 36000 NUMBERS



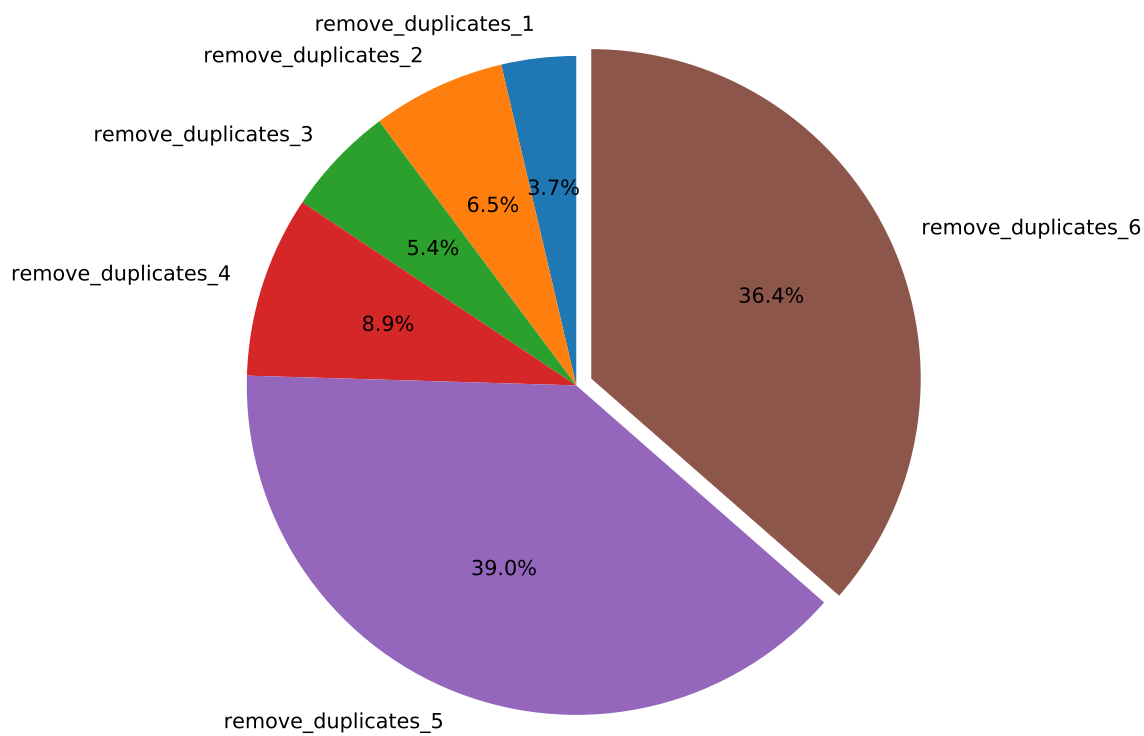
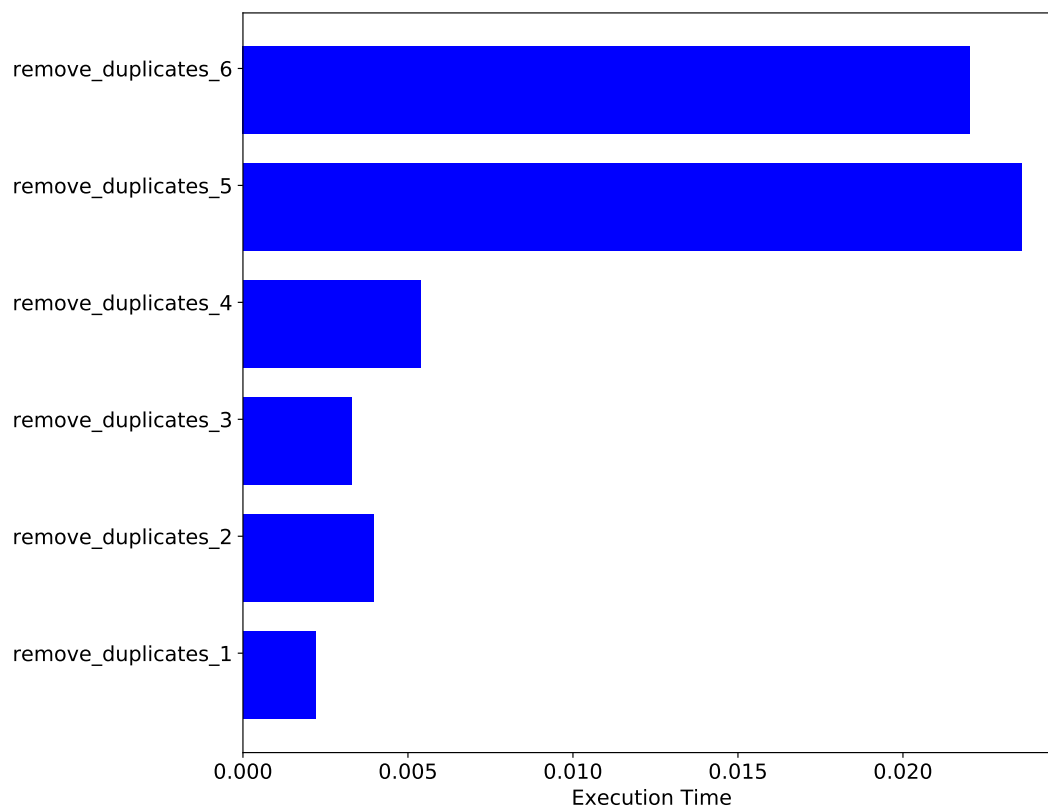
TIMES FOR THE FIRST 38000 NUMBERS



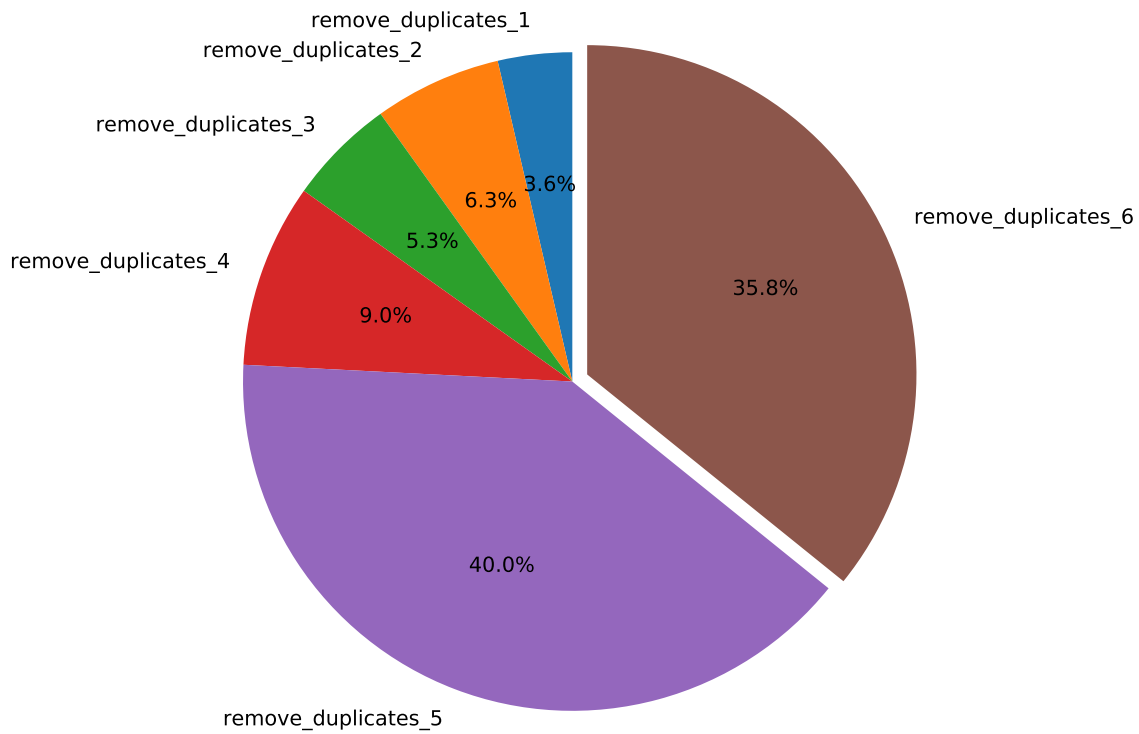
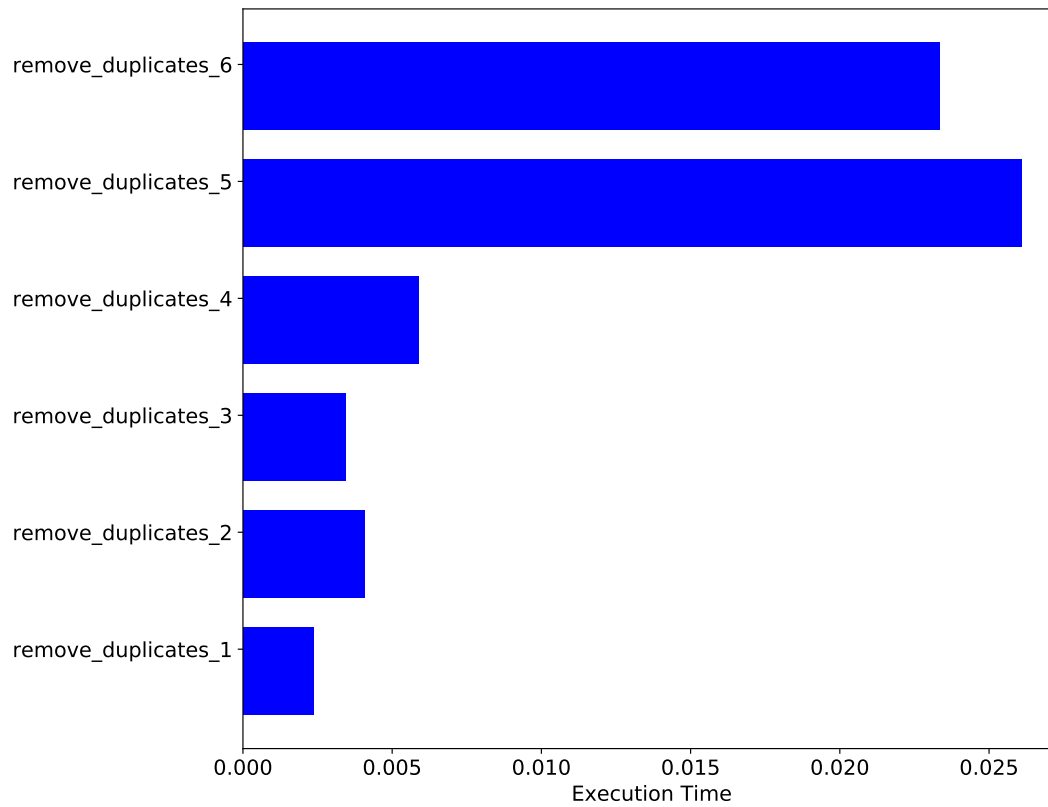
TIMES FOR THE FIRST 40000 NUMBERS



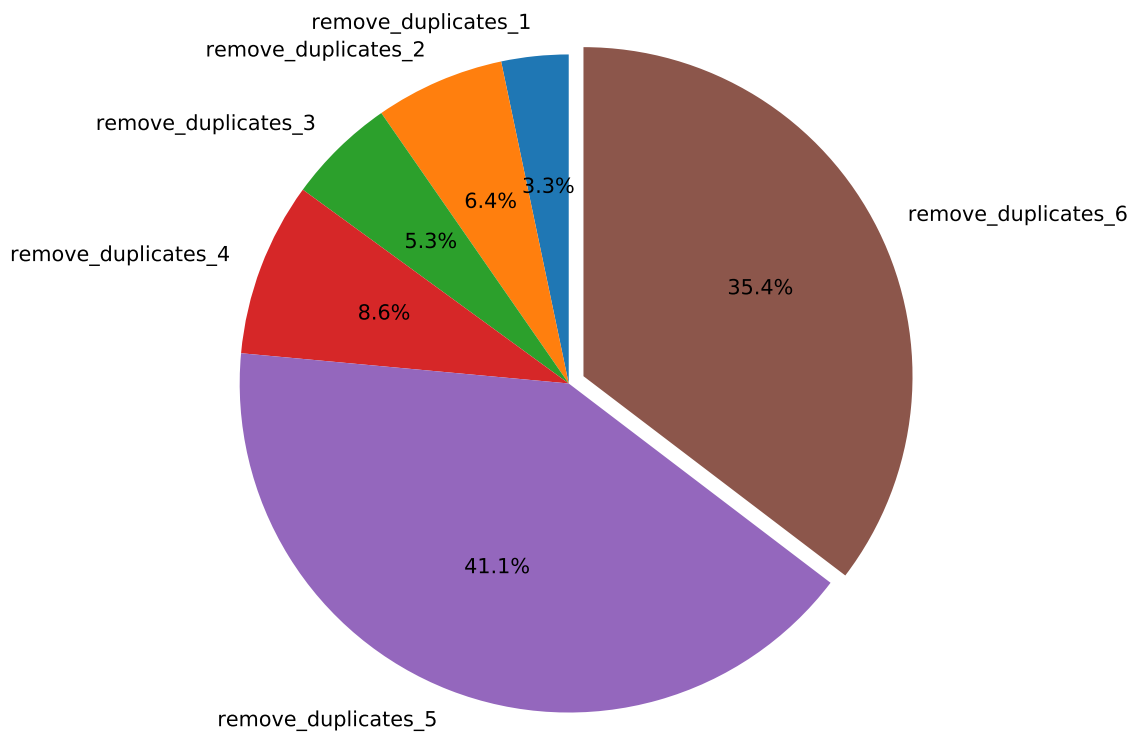
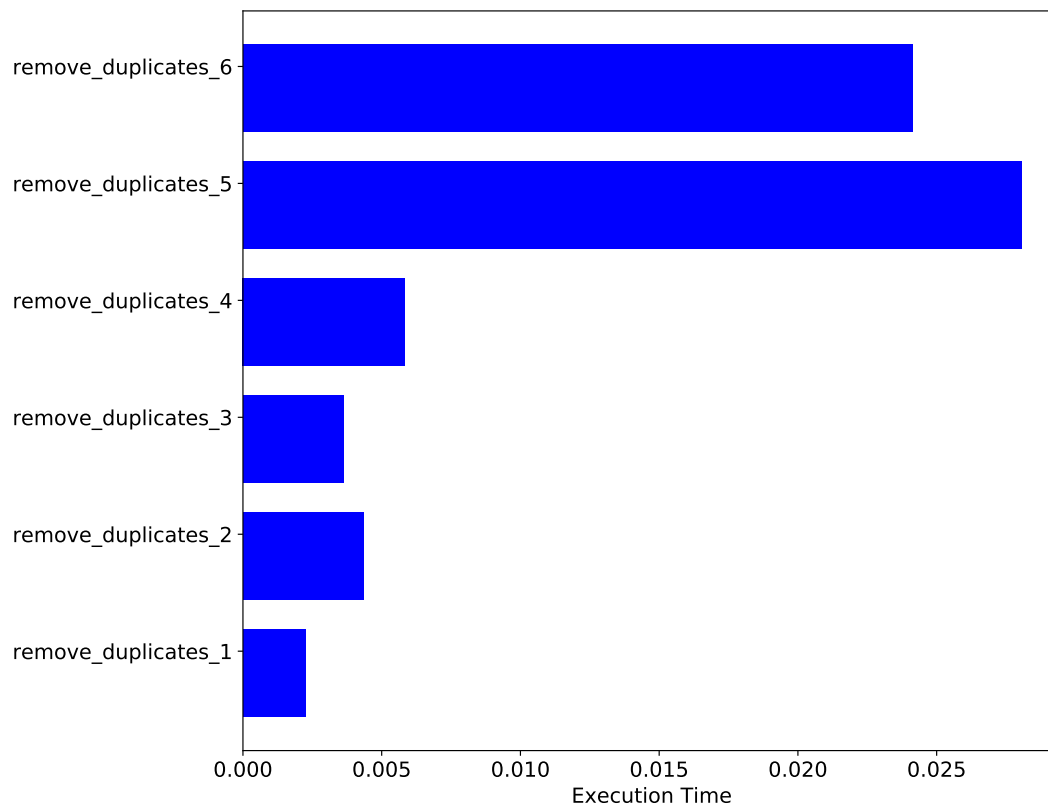
TIMES FOR THE FIRST 42000 NUMBERS



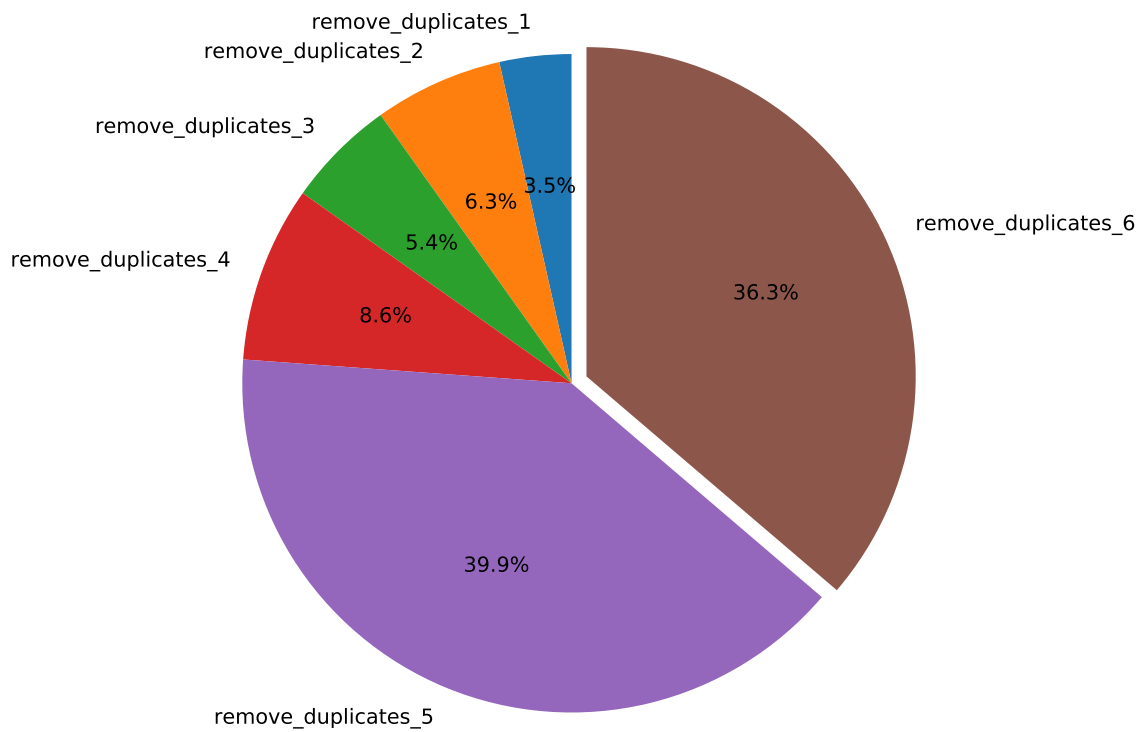
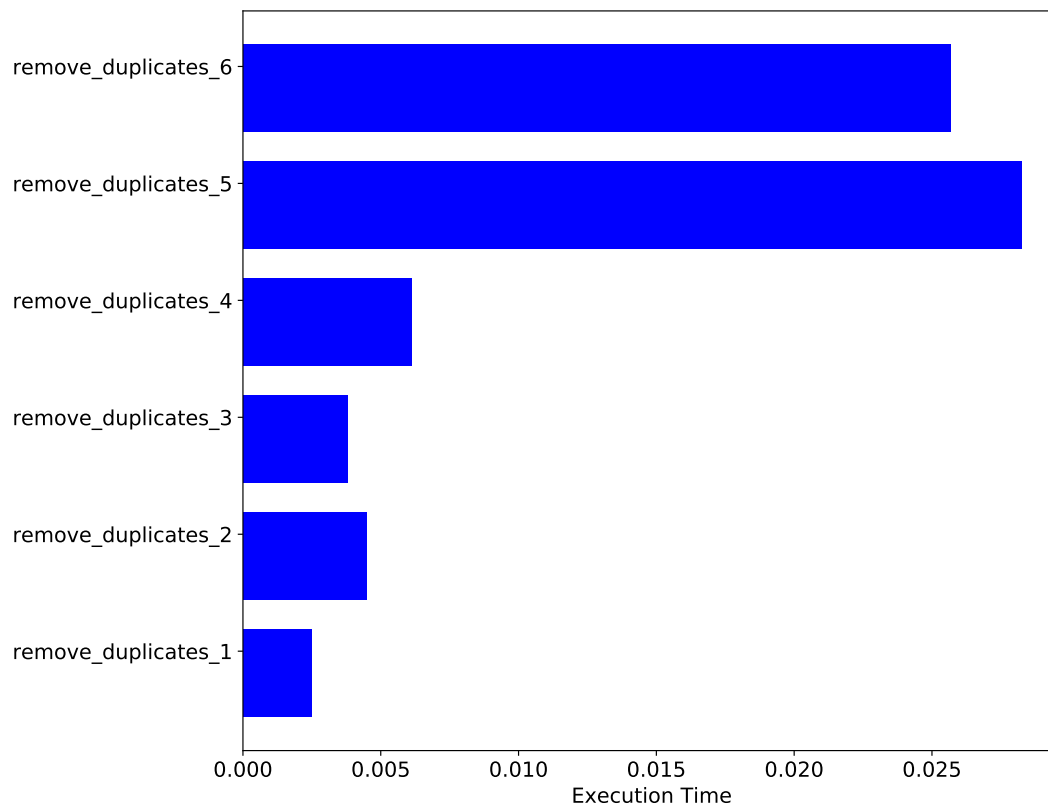
TIMES FOR THE FIRST 44000 NUMBERS



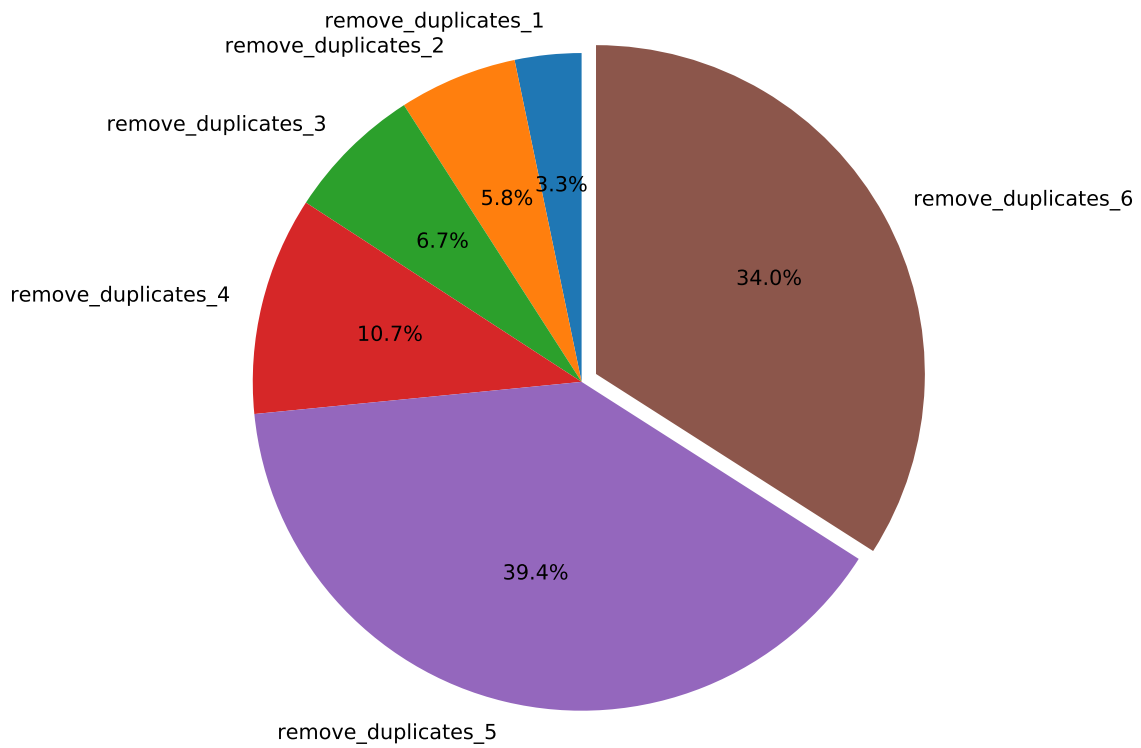
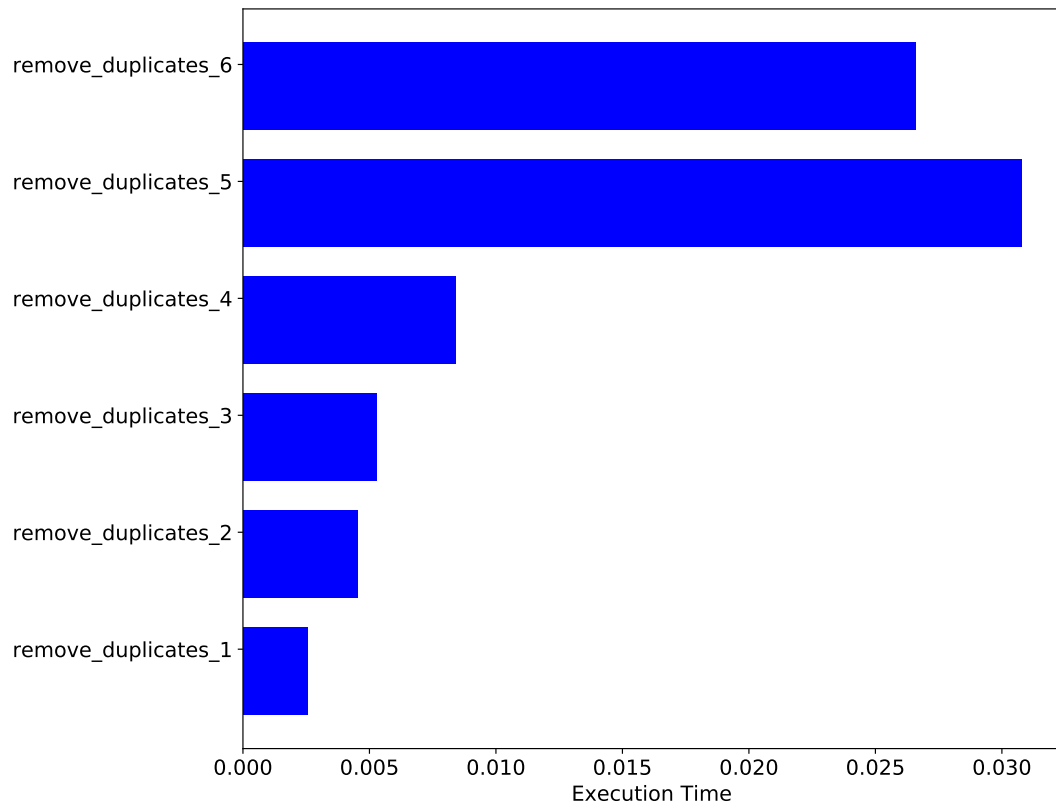
TIMES FOR THE FIRST 46000 NUMBERS



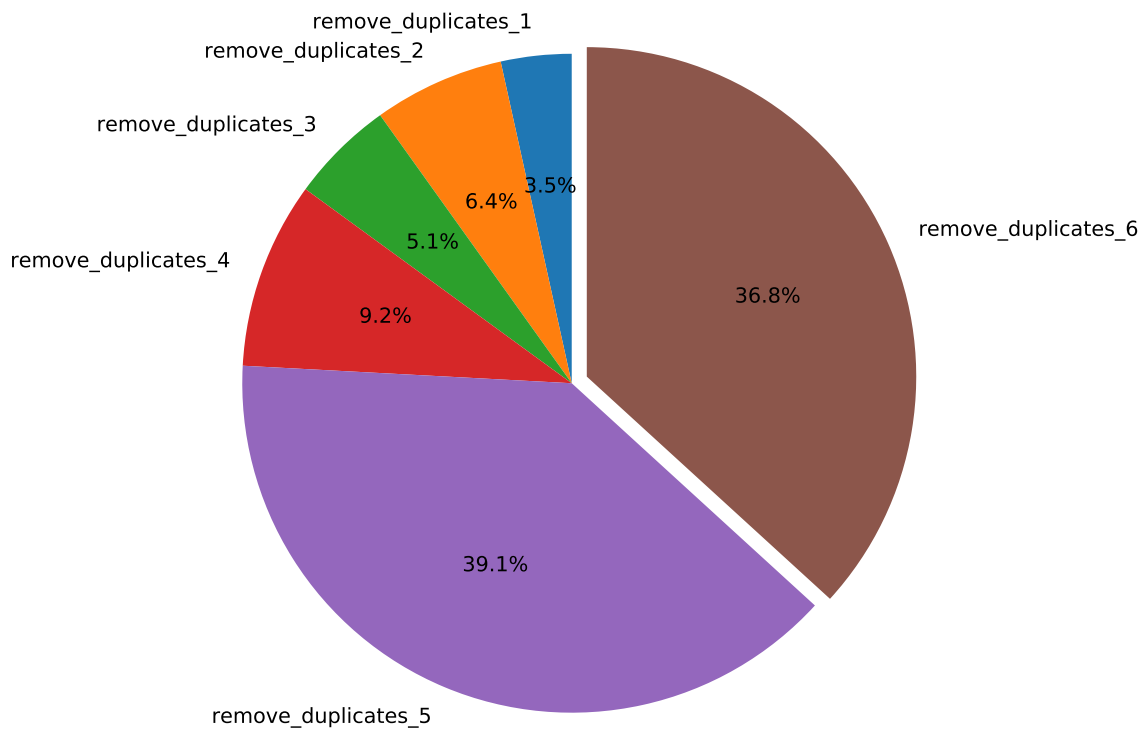
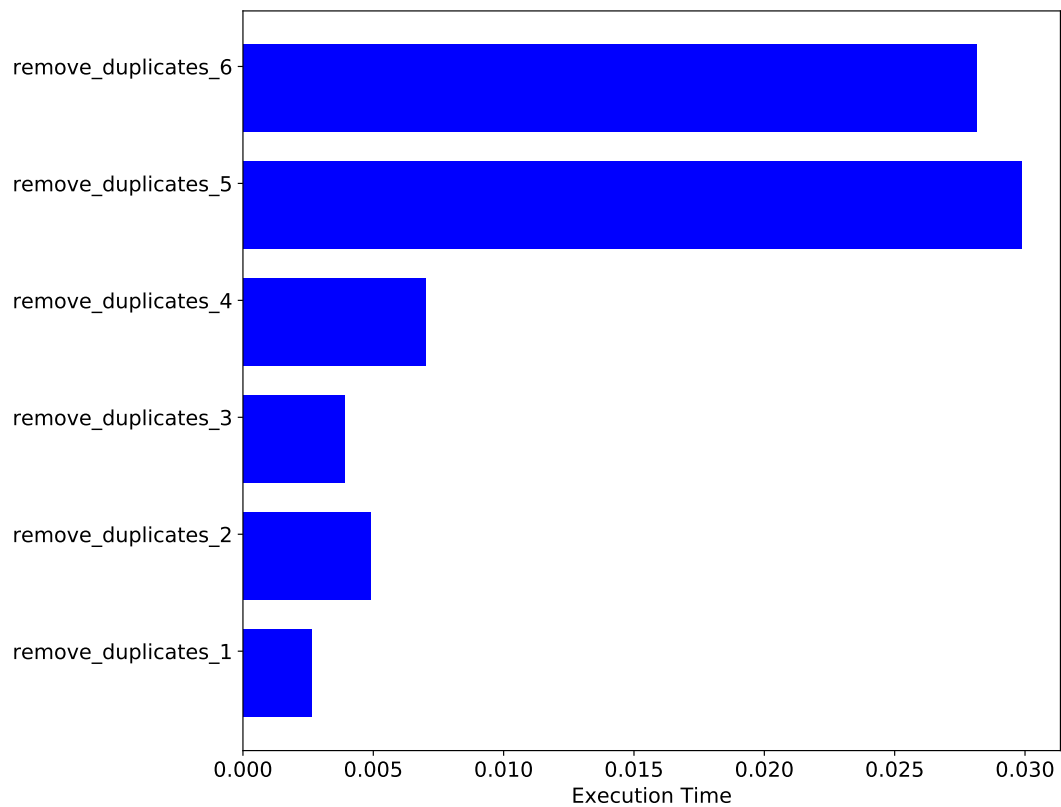
TIMES FOR THE FIRST 48000 NUMBERS



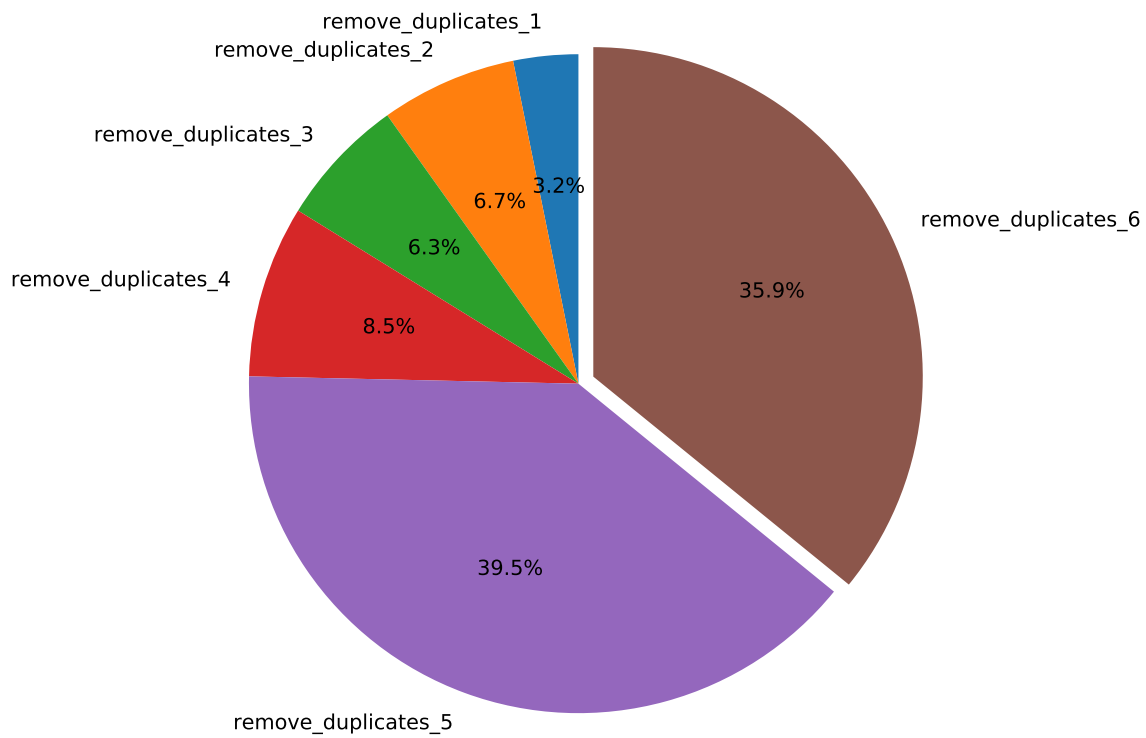
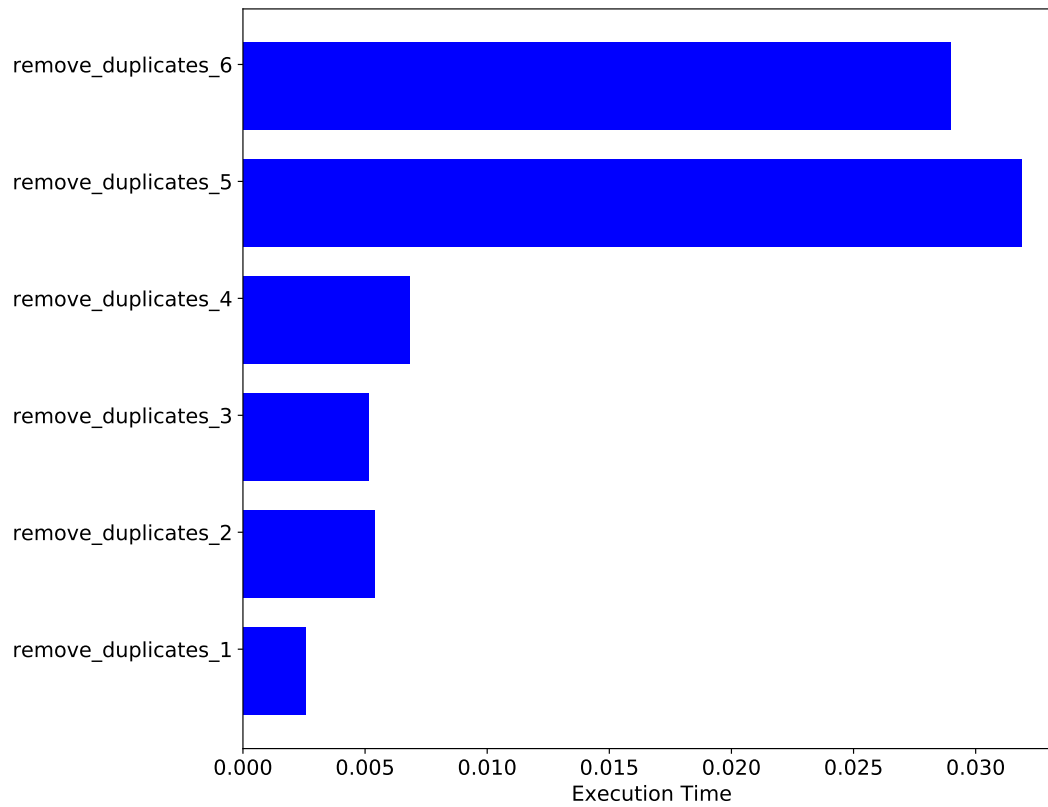
TIMES FOR THE FIRST 50000 NUMBERS



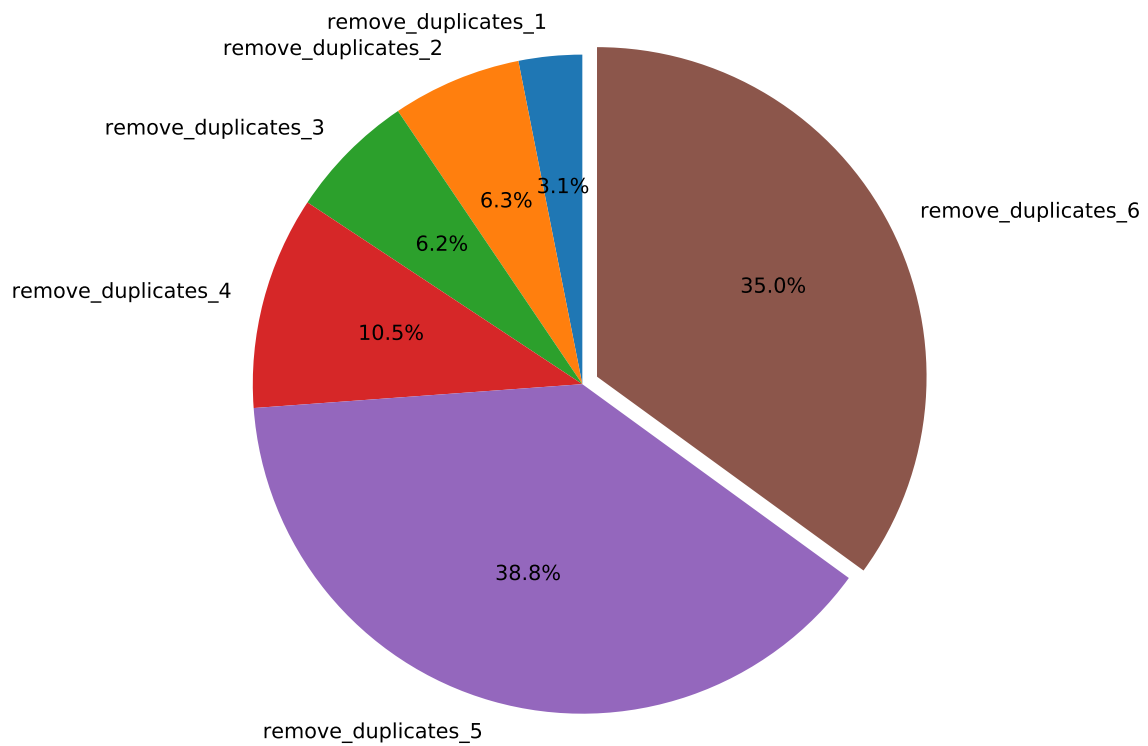
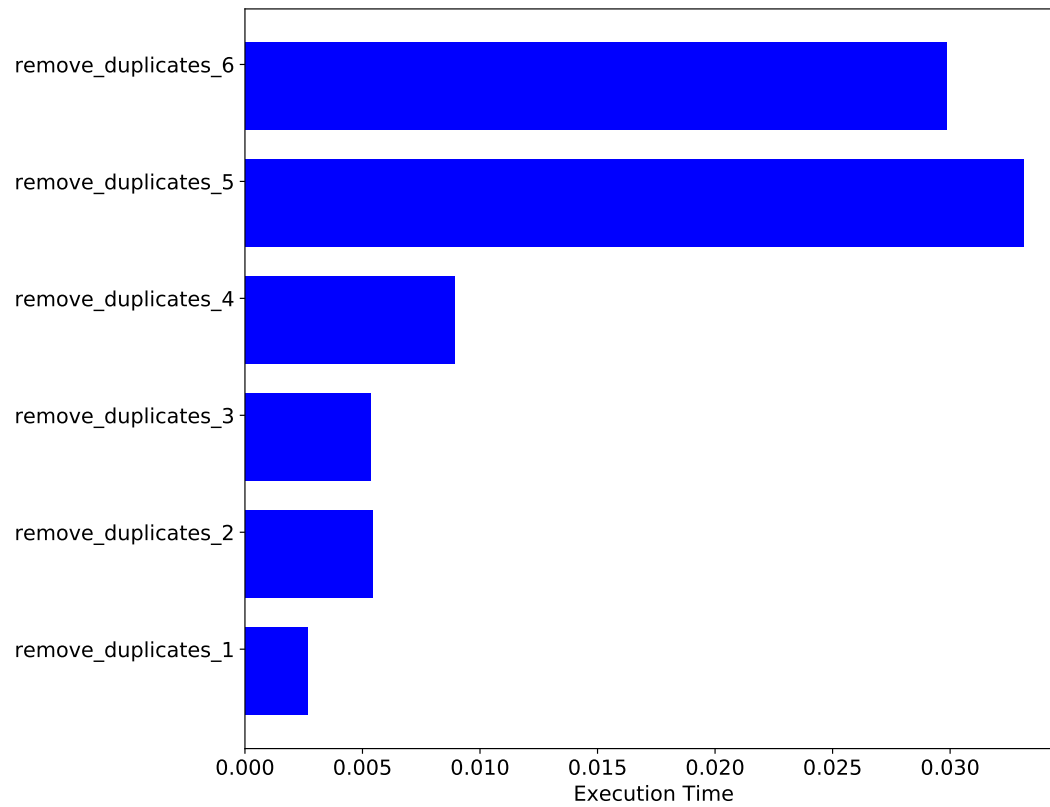
TIMES FOR THE FIRST 52000 NUMBERS



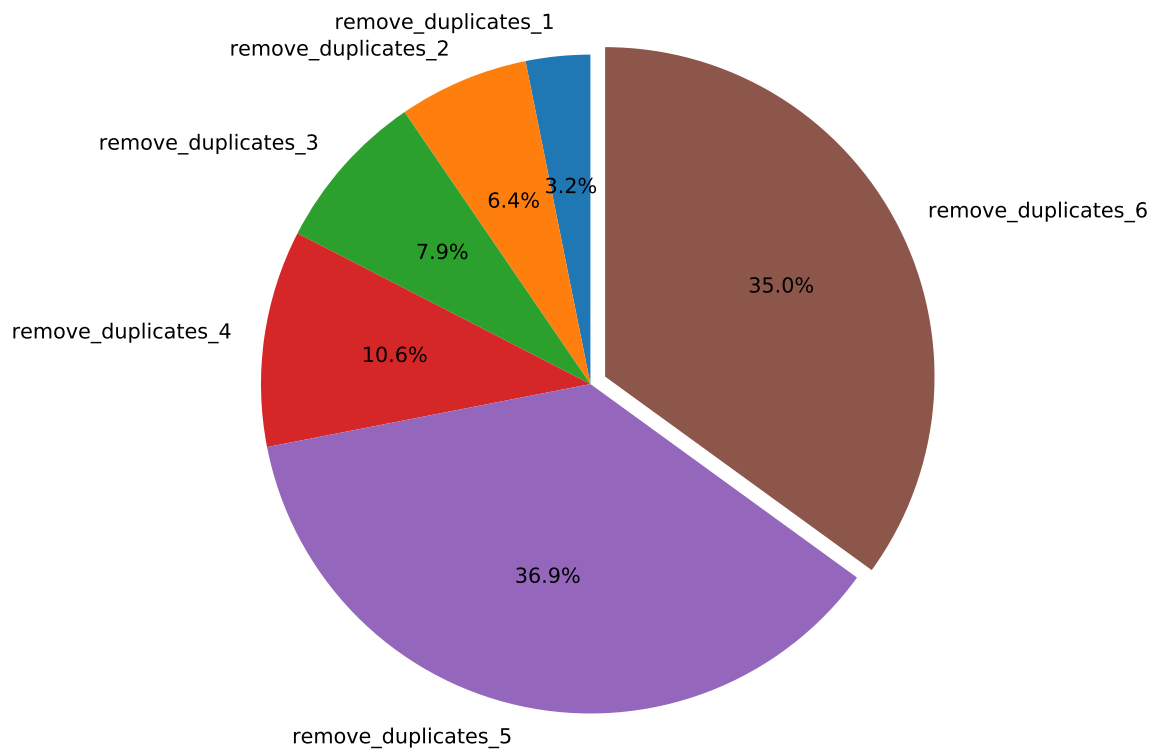
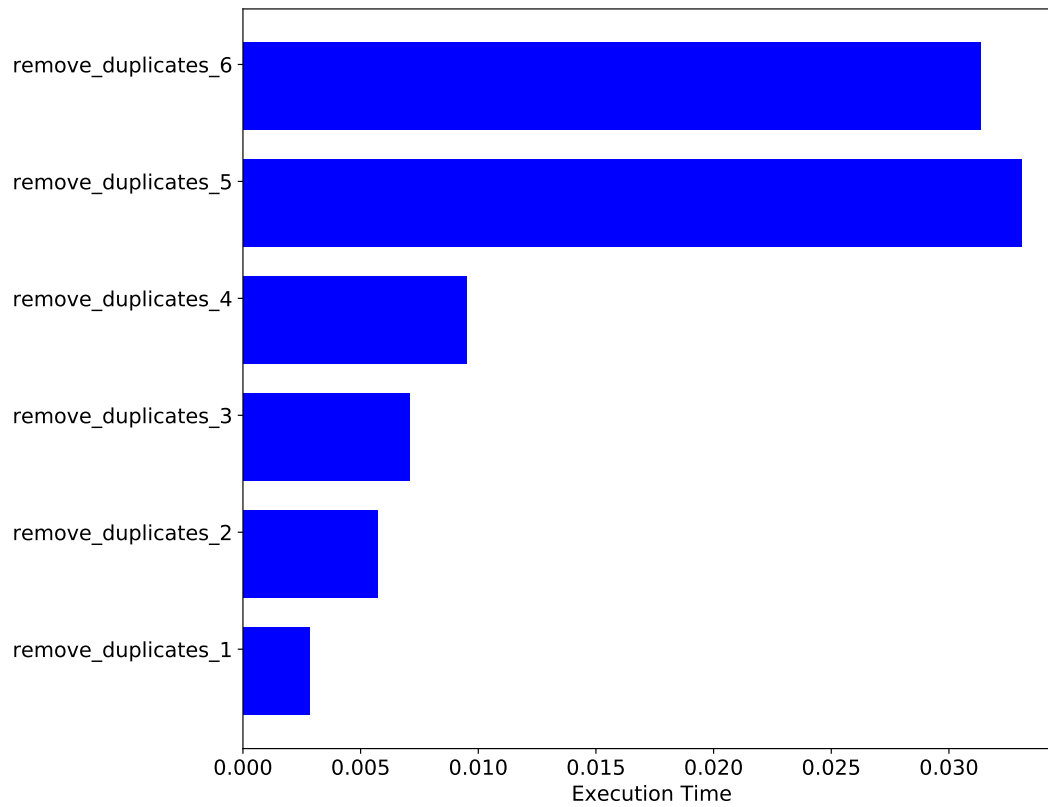
TIMES FOR THE FIRST 54000 NUMBERS



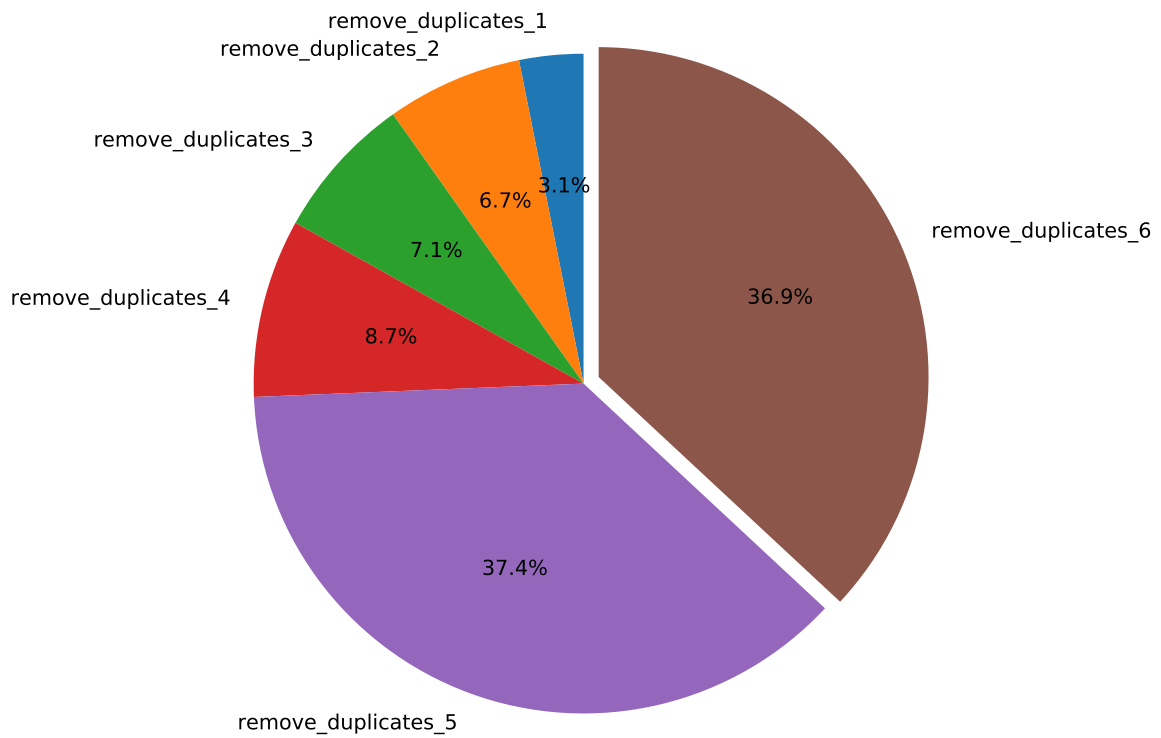
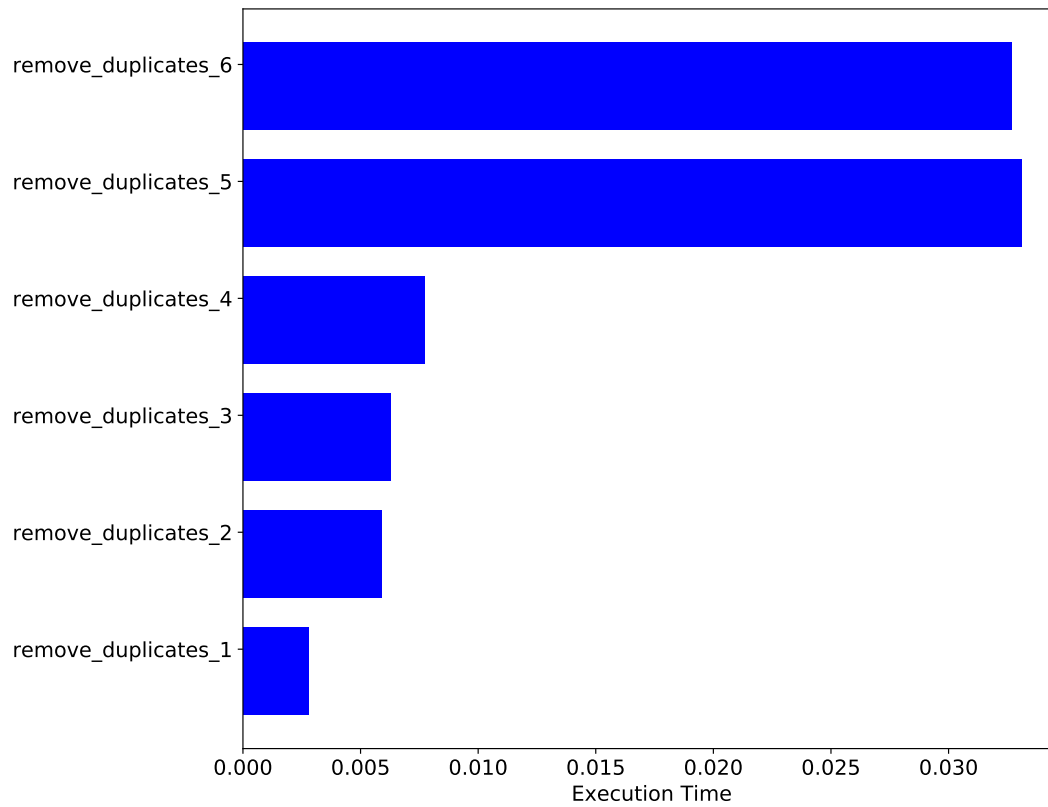
TIMES FOR THE FIRST 56000 NUMBERS



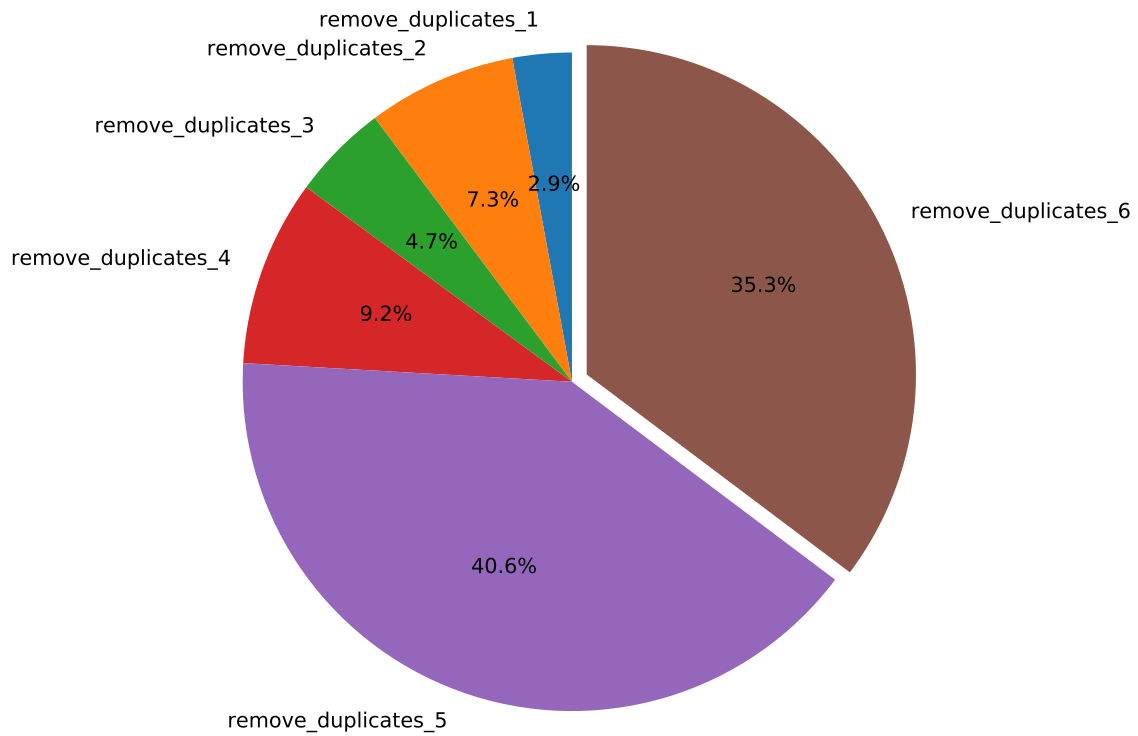
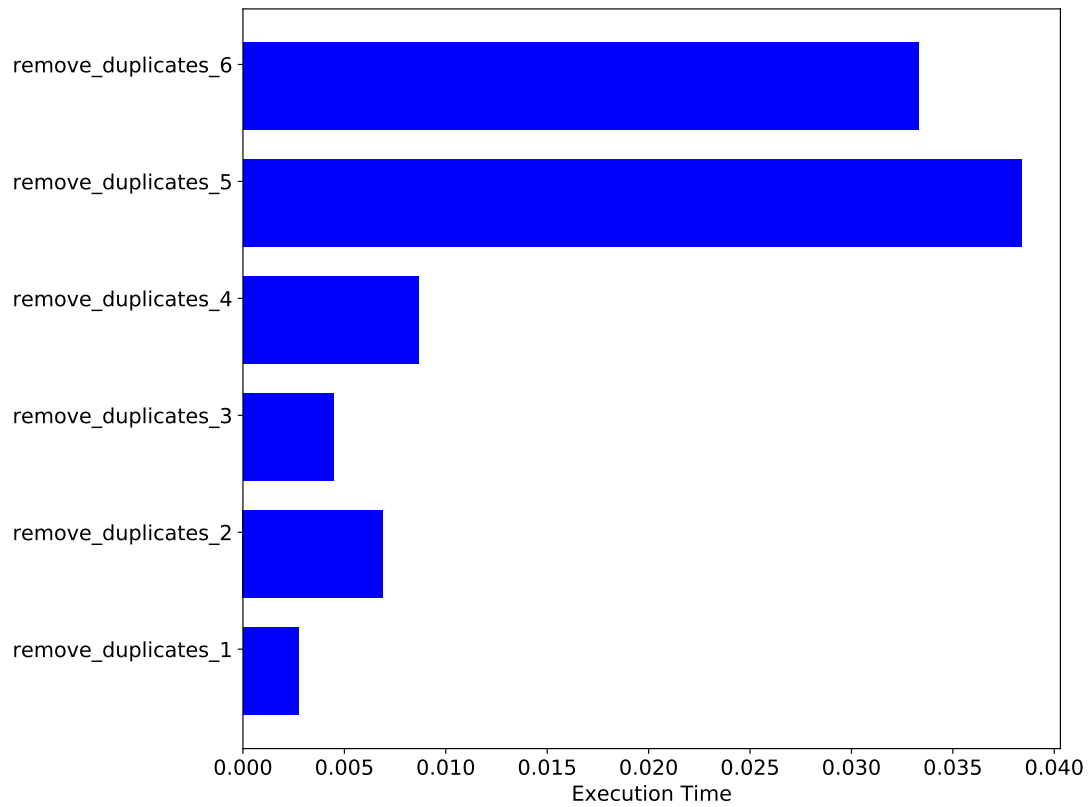
TIMES FOR THE FIRST 58000 NUMBERS



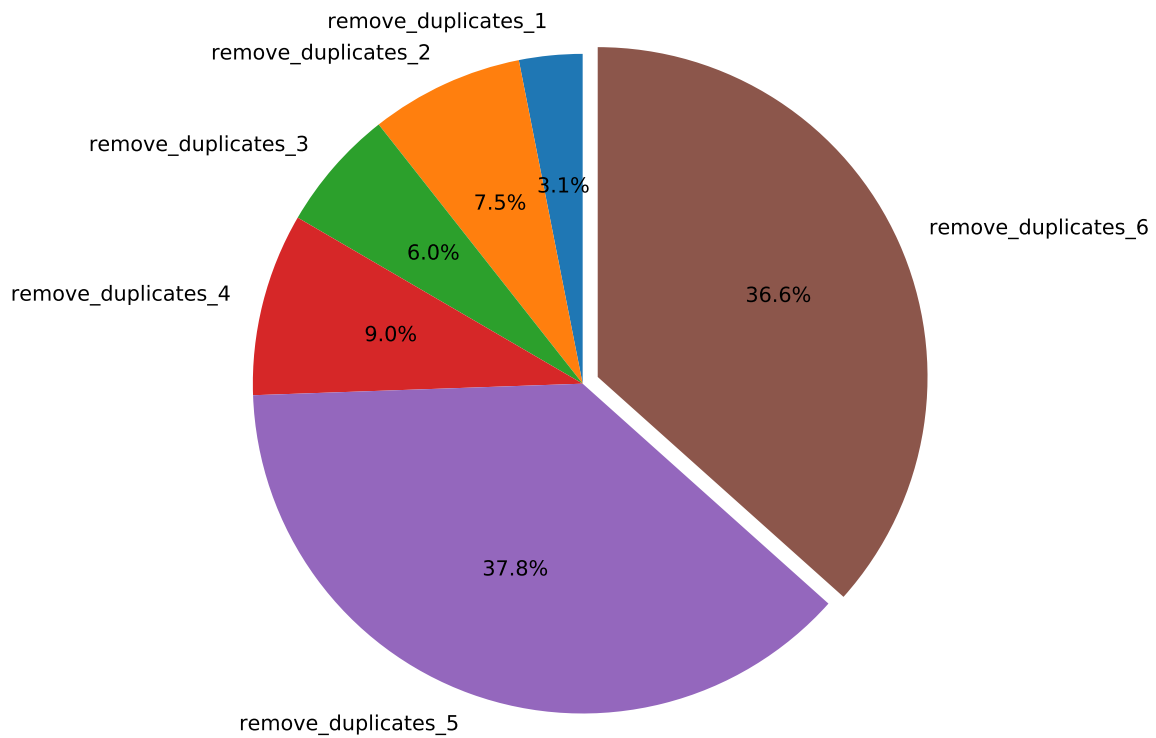
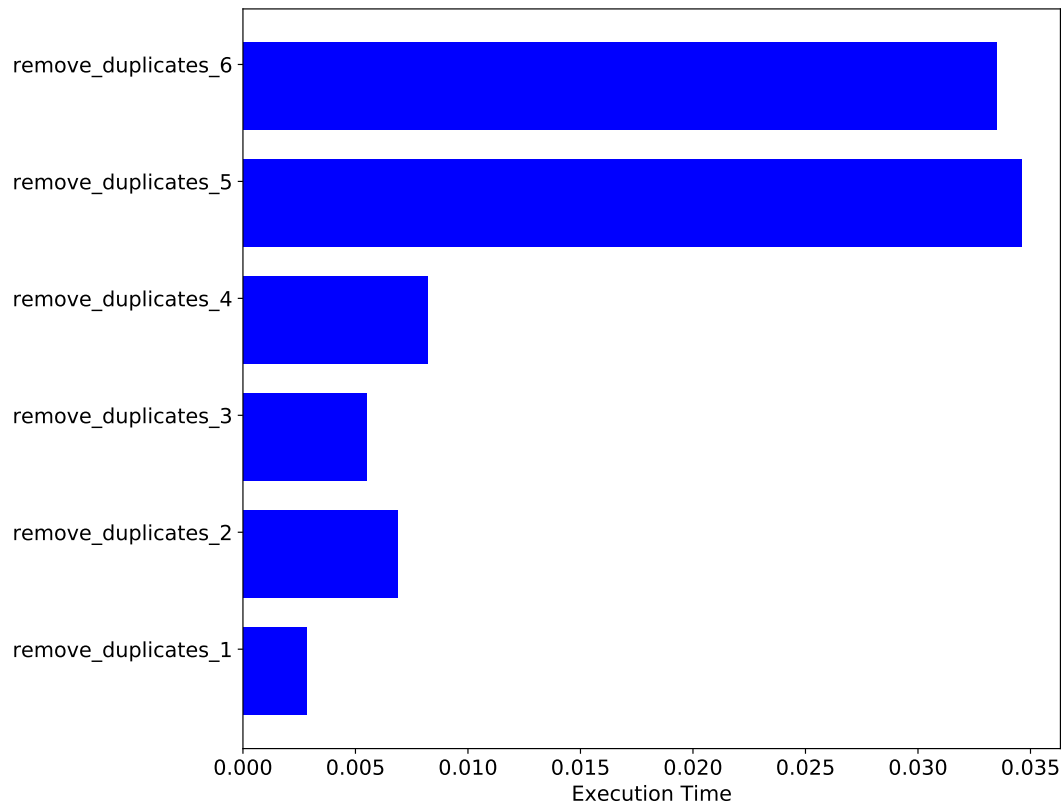
TIMES FOR THE FIRST 60000 NUMBERS



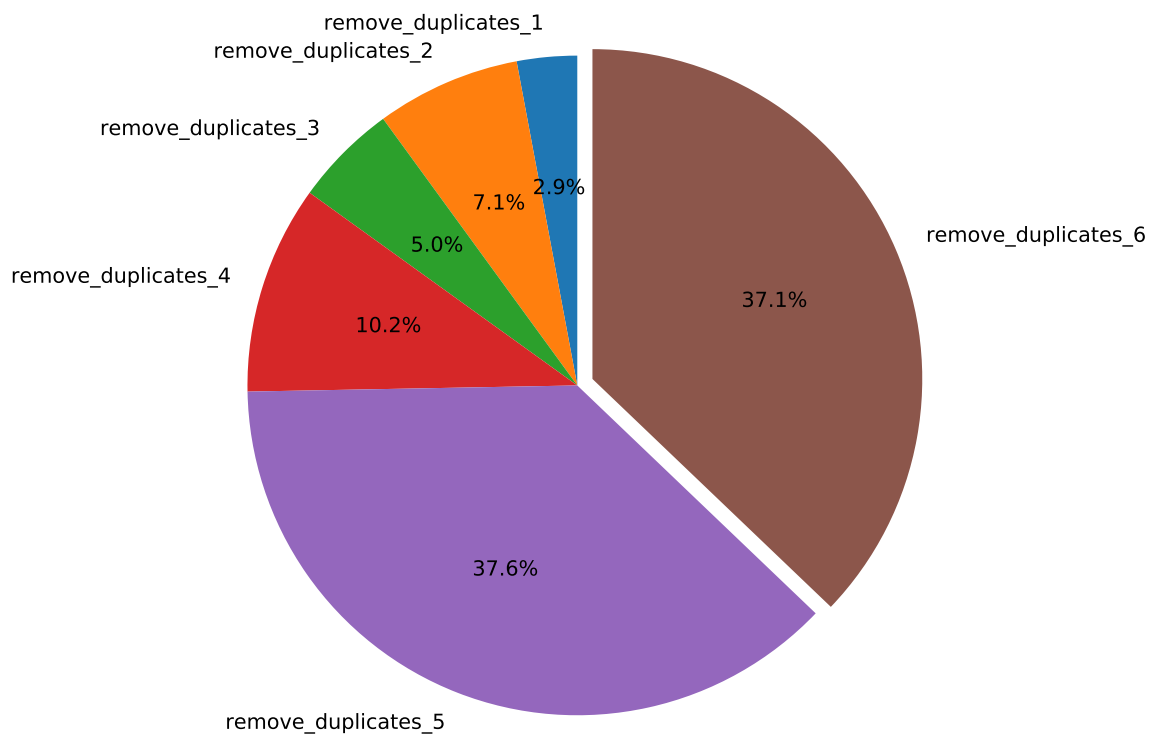
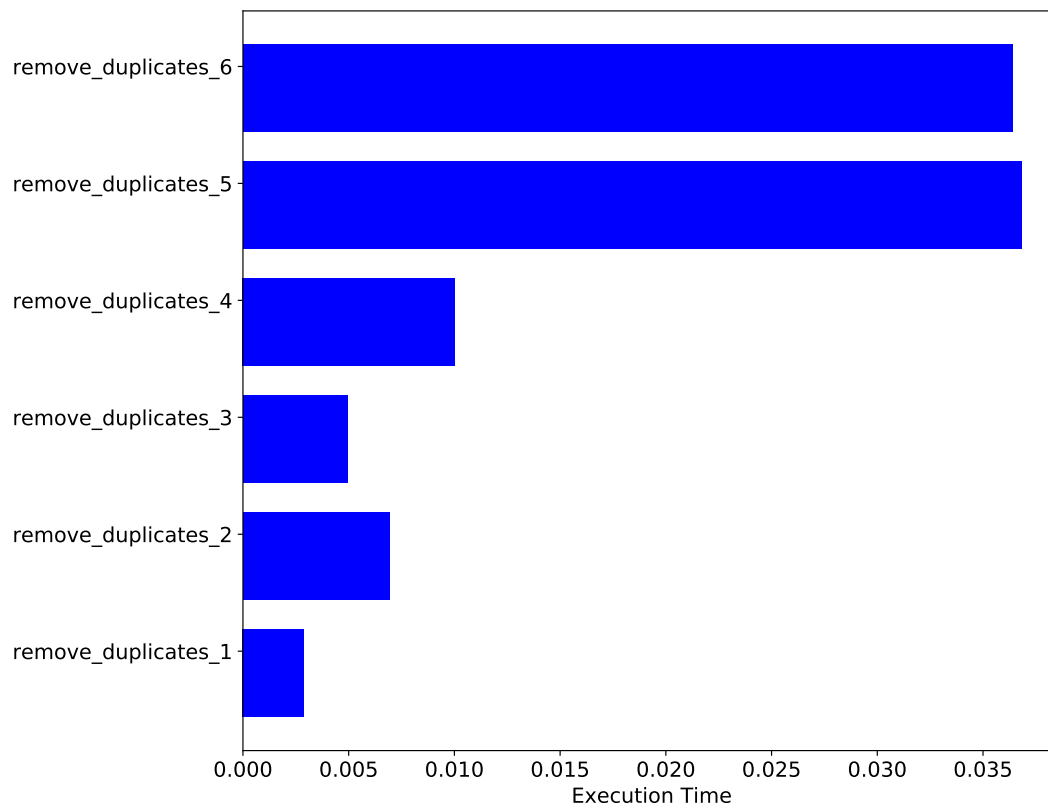
TIMES FOR THE FIRST 62000 NUMBERS



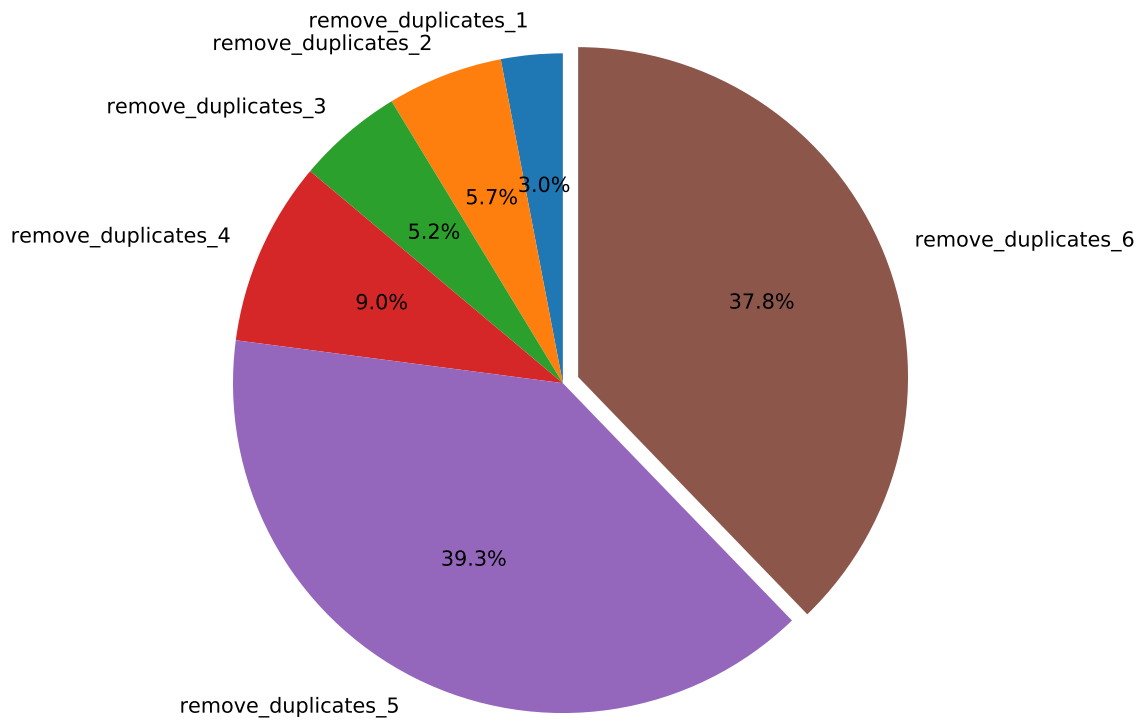
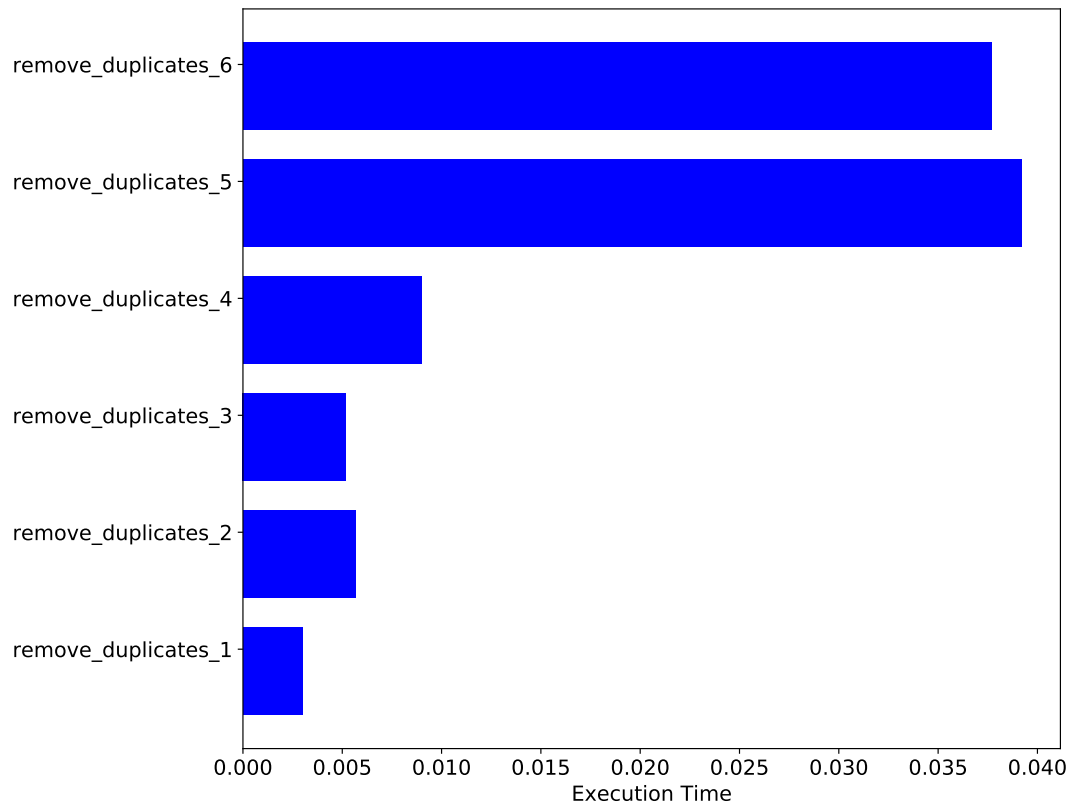
TIMES FOR THE FIRST 64000 NUMBERS



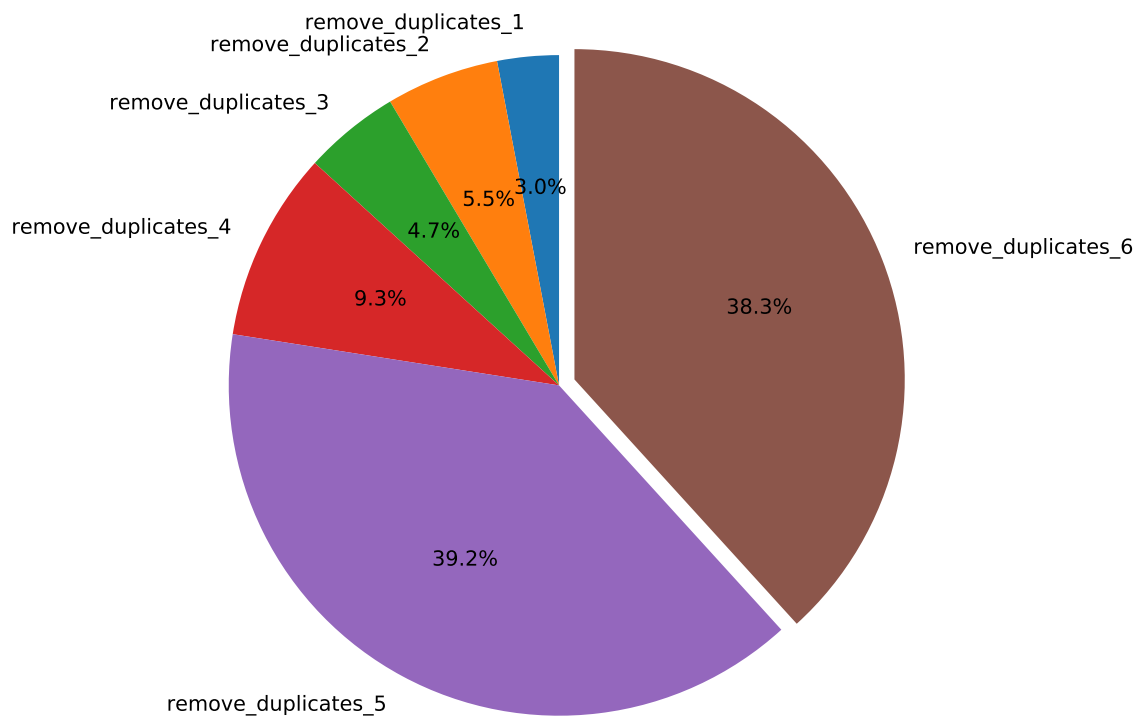
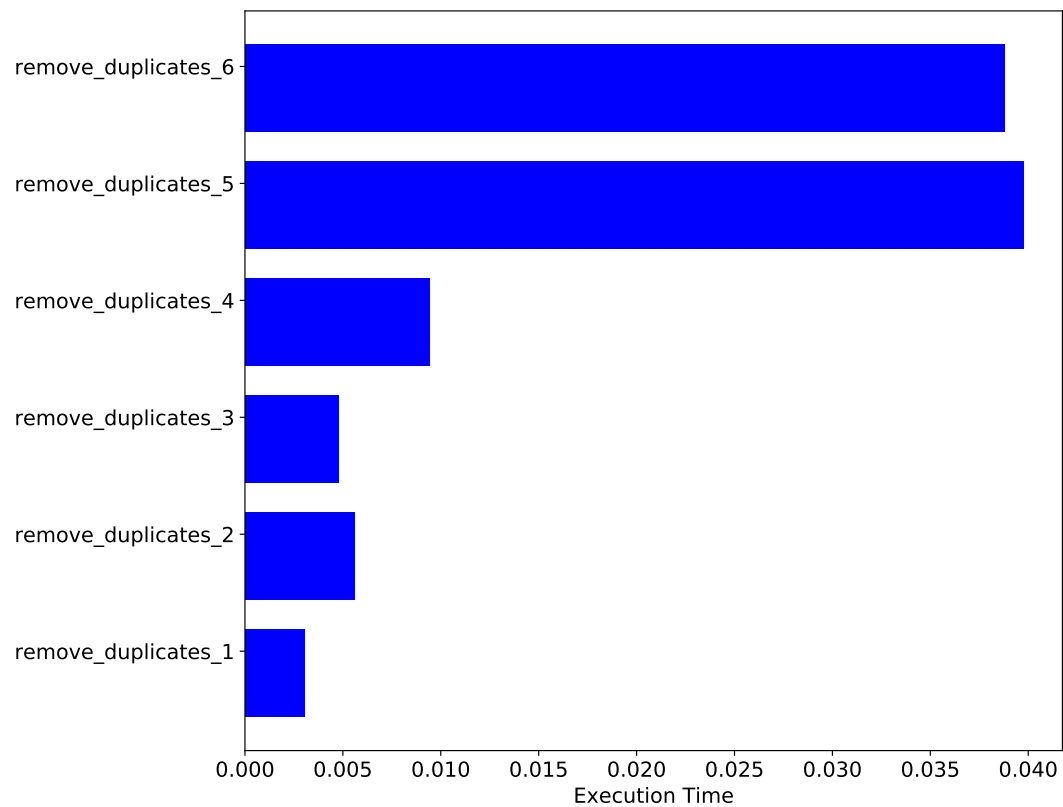
TIMES FOR THE FIRST 66000 NUMBERS



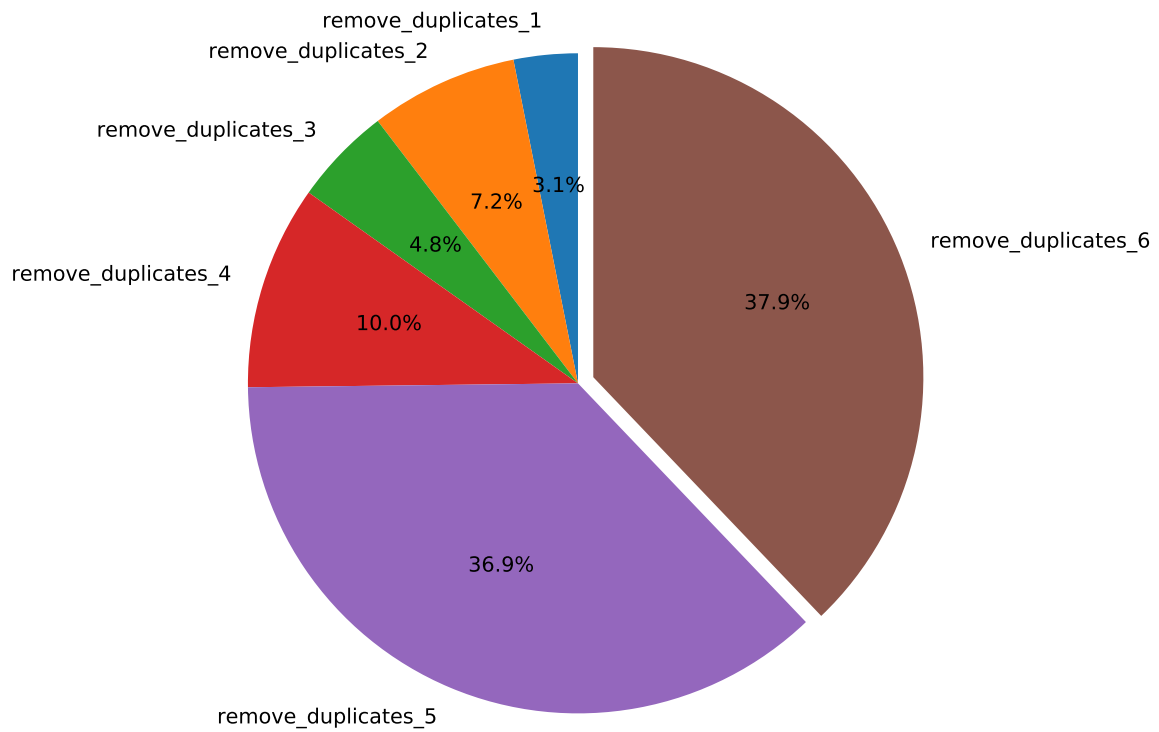
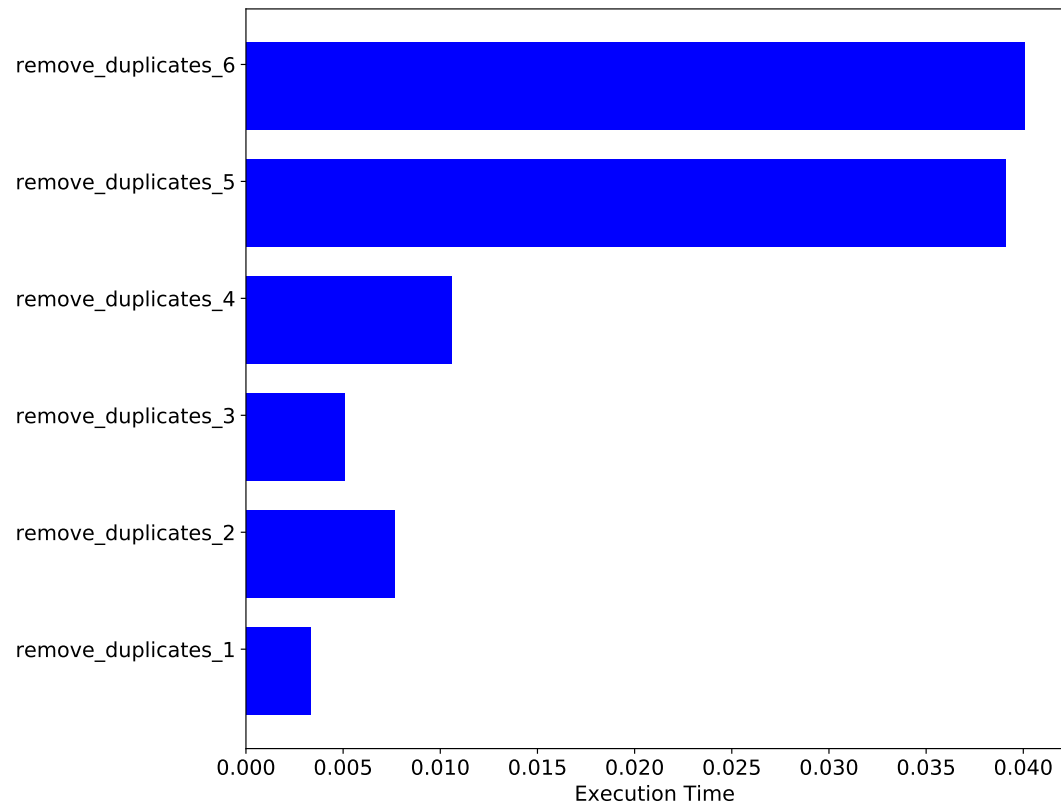
TIMES FOR THE FIRST 68000 NUMBERS



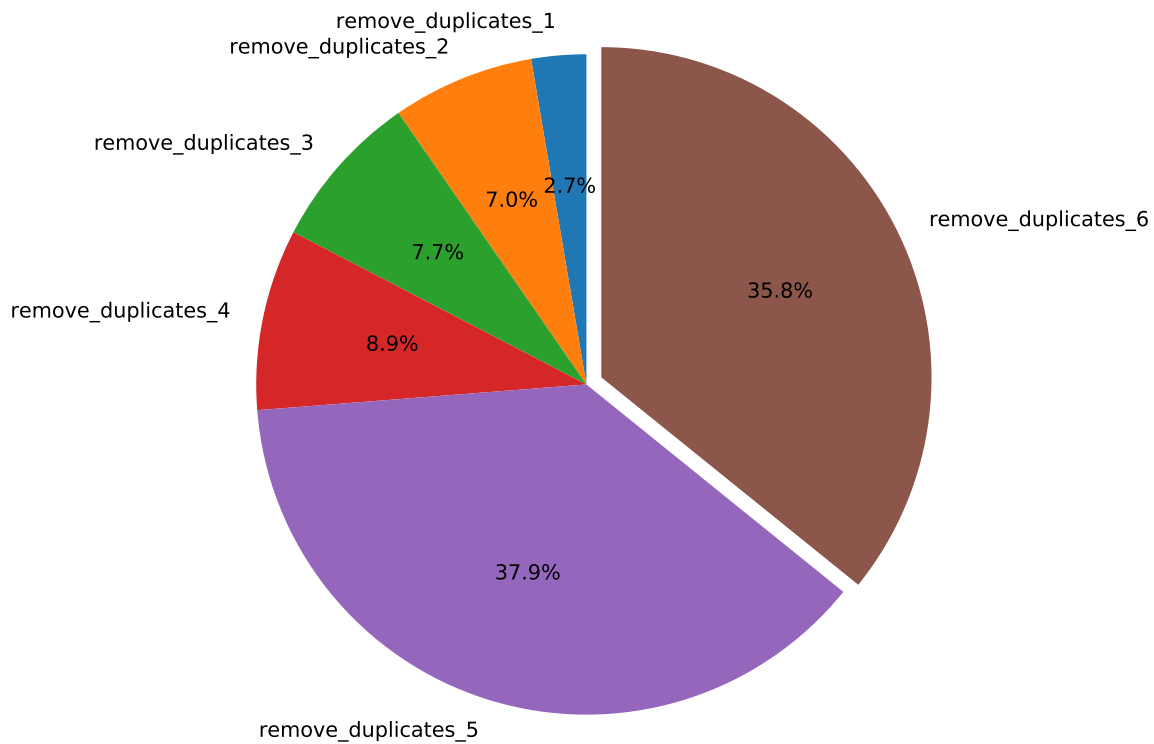
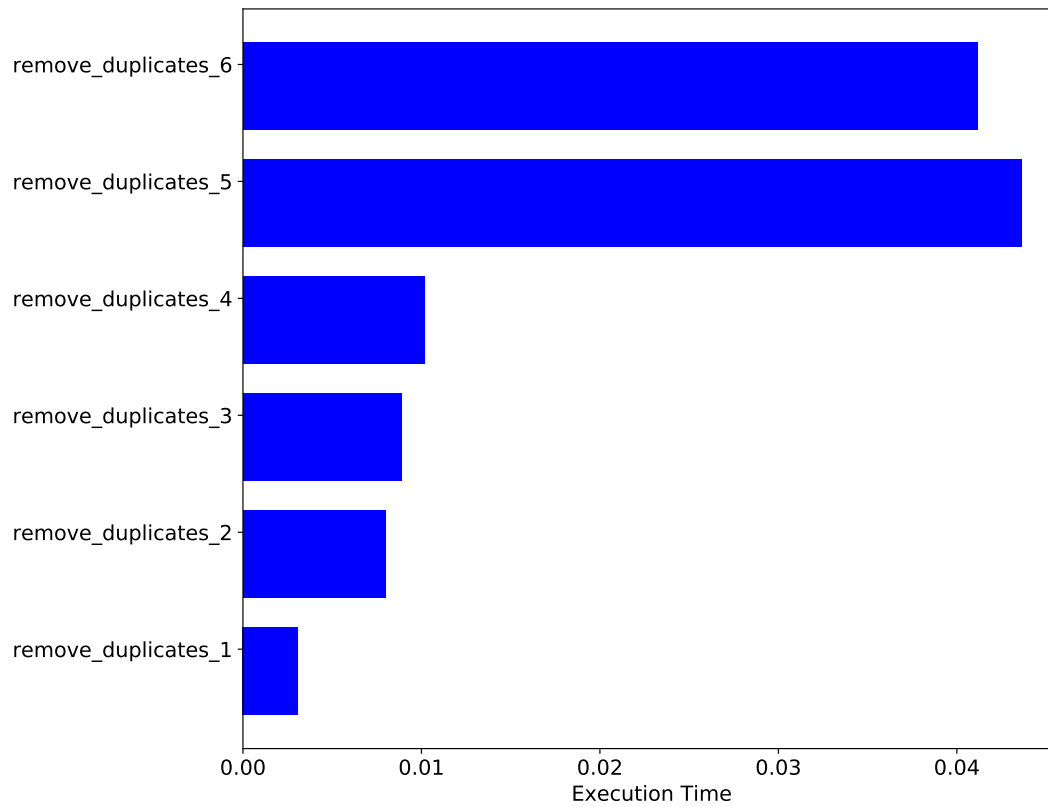
TIMES FOR THE FIRST 70000 NUMBERS



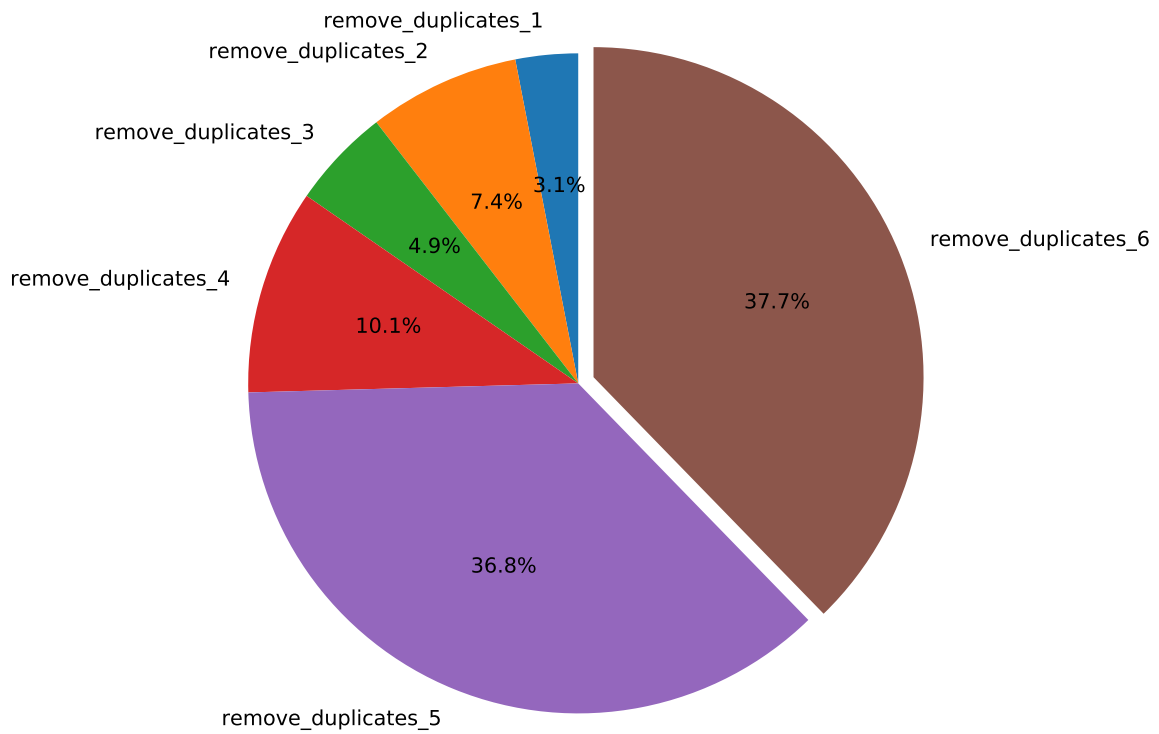
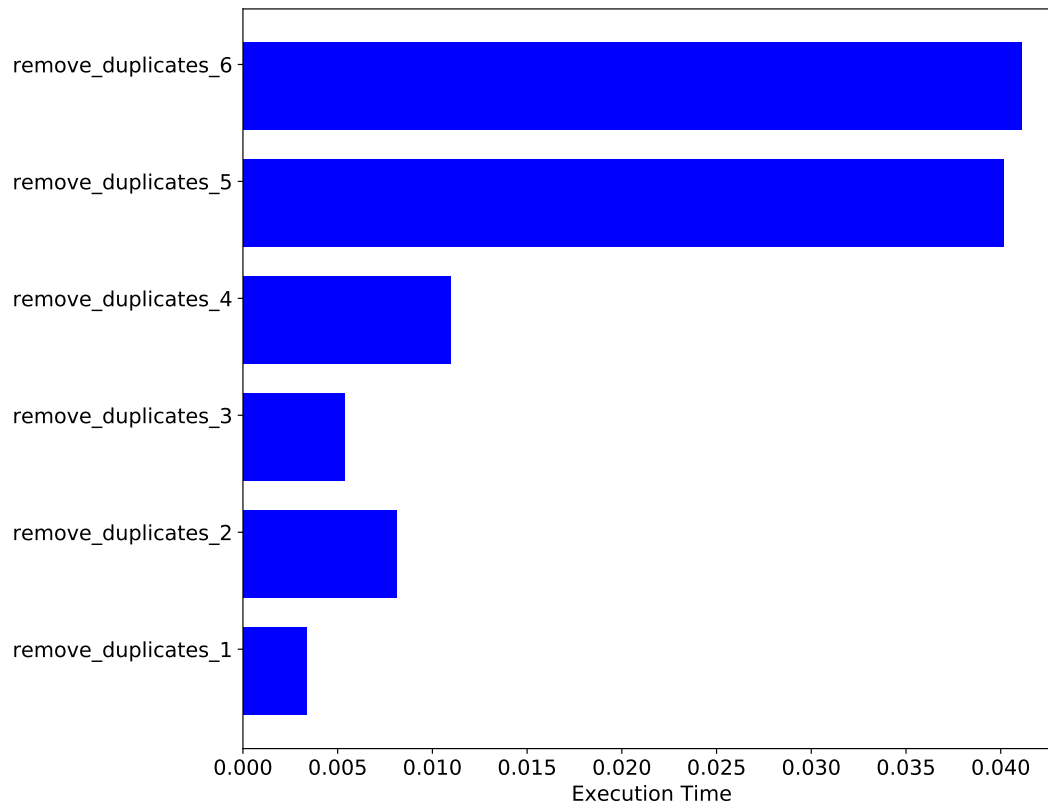
TIMES FOR THE FIRST 72000 NUMBERS



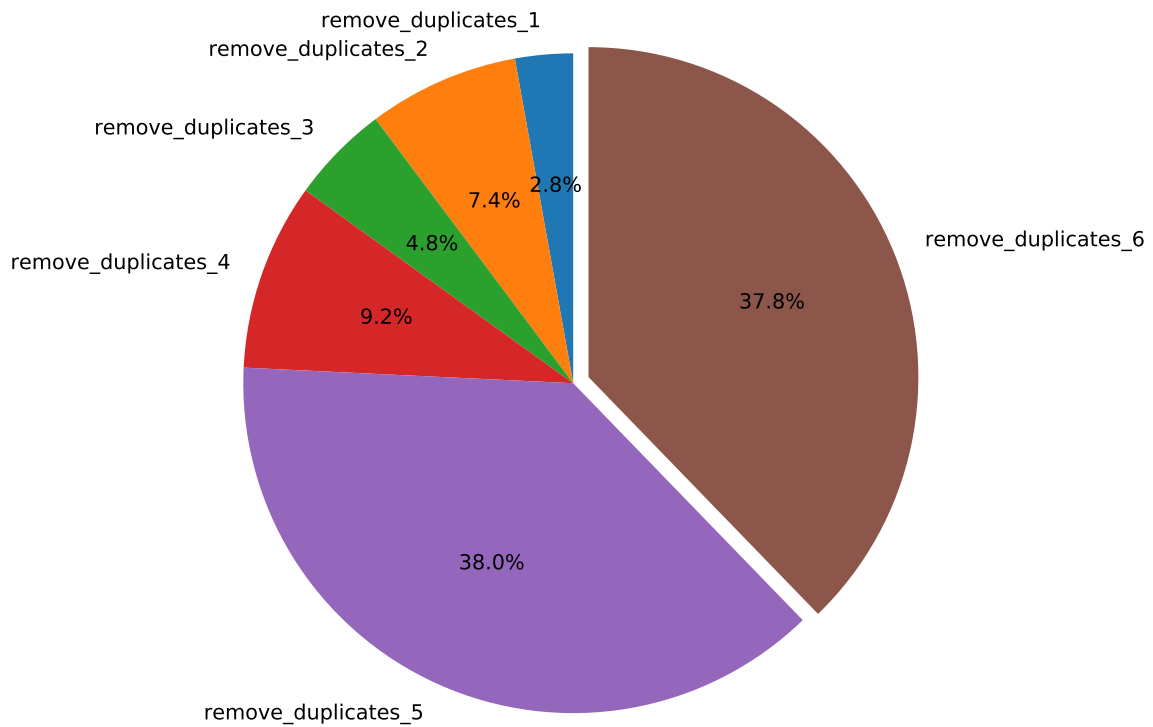
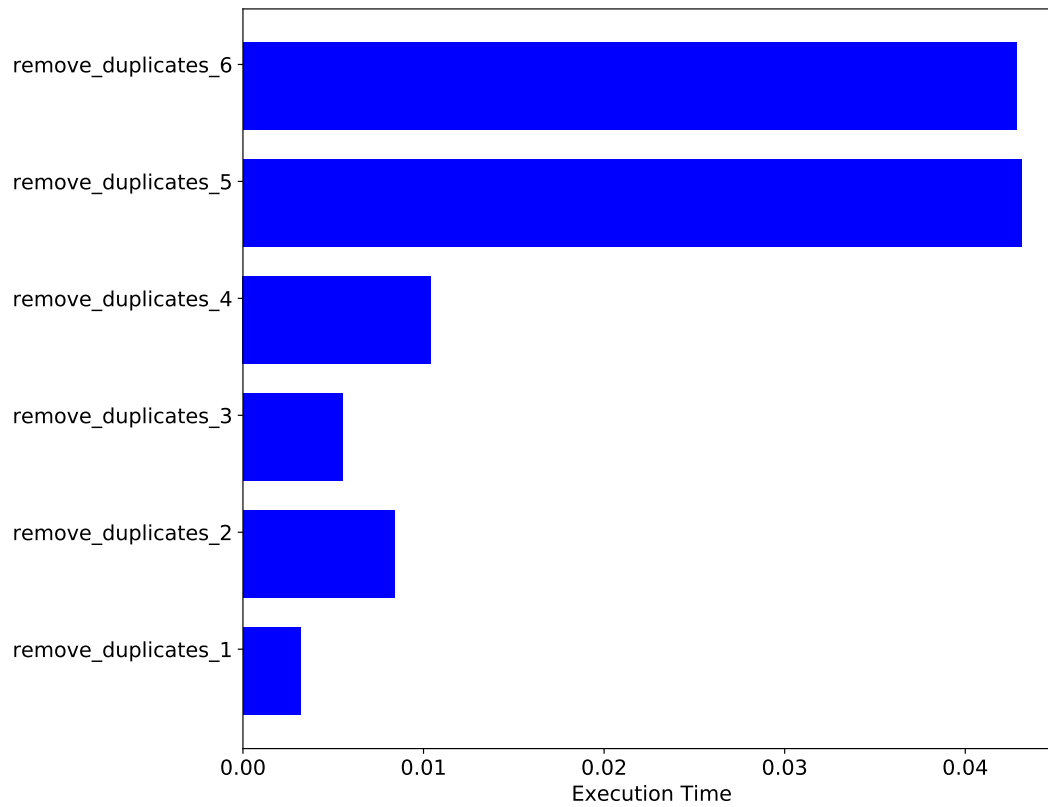
TIMES FOR THE FIRST 74000 NUMBERS



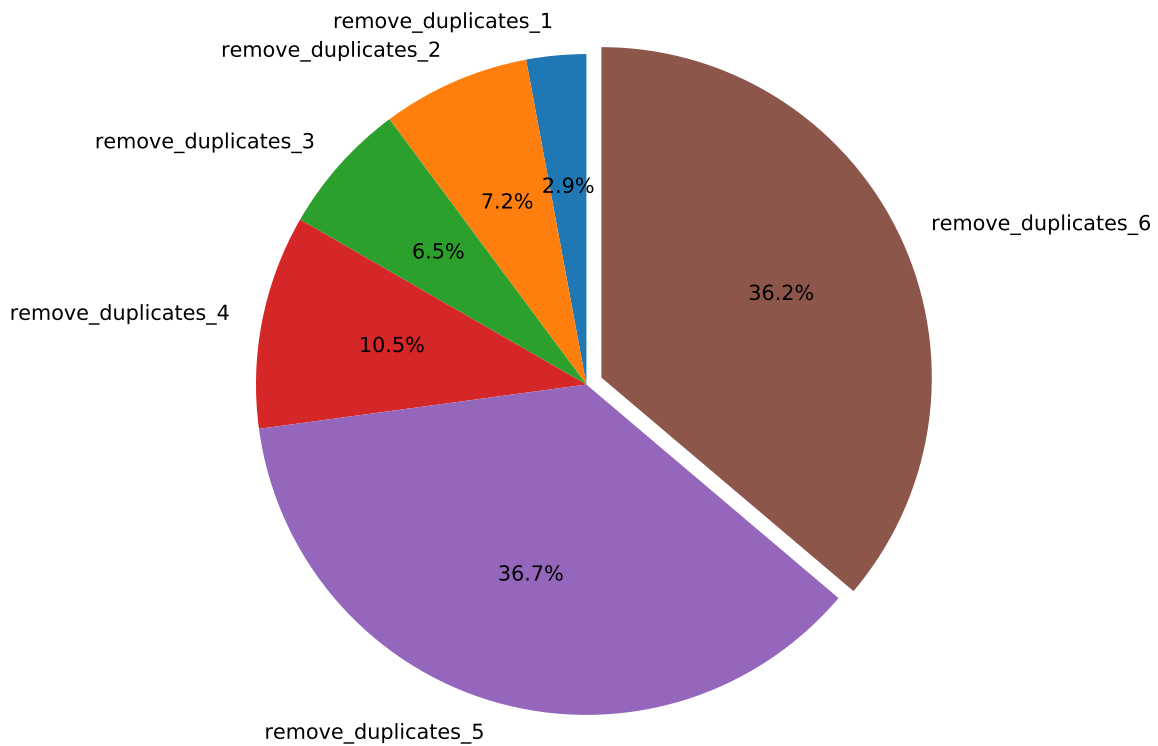
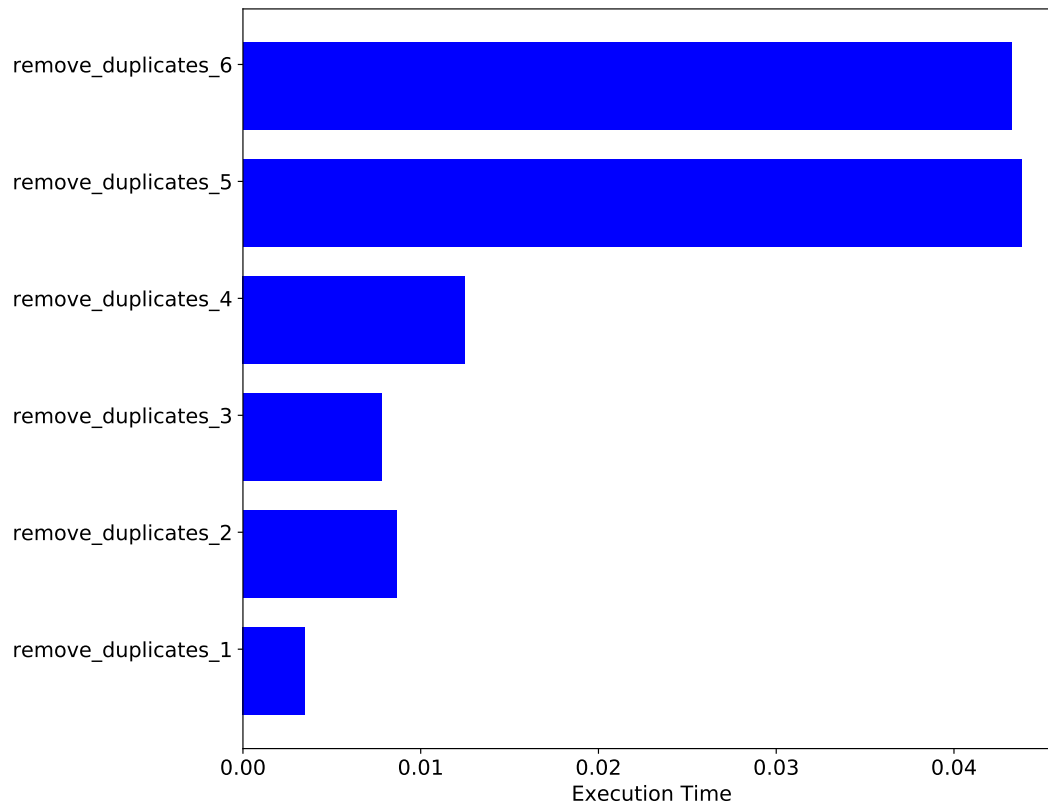
TIMES FOR THE FIRST 76000 NUMBERS



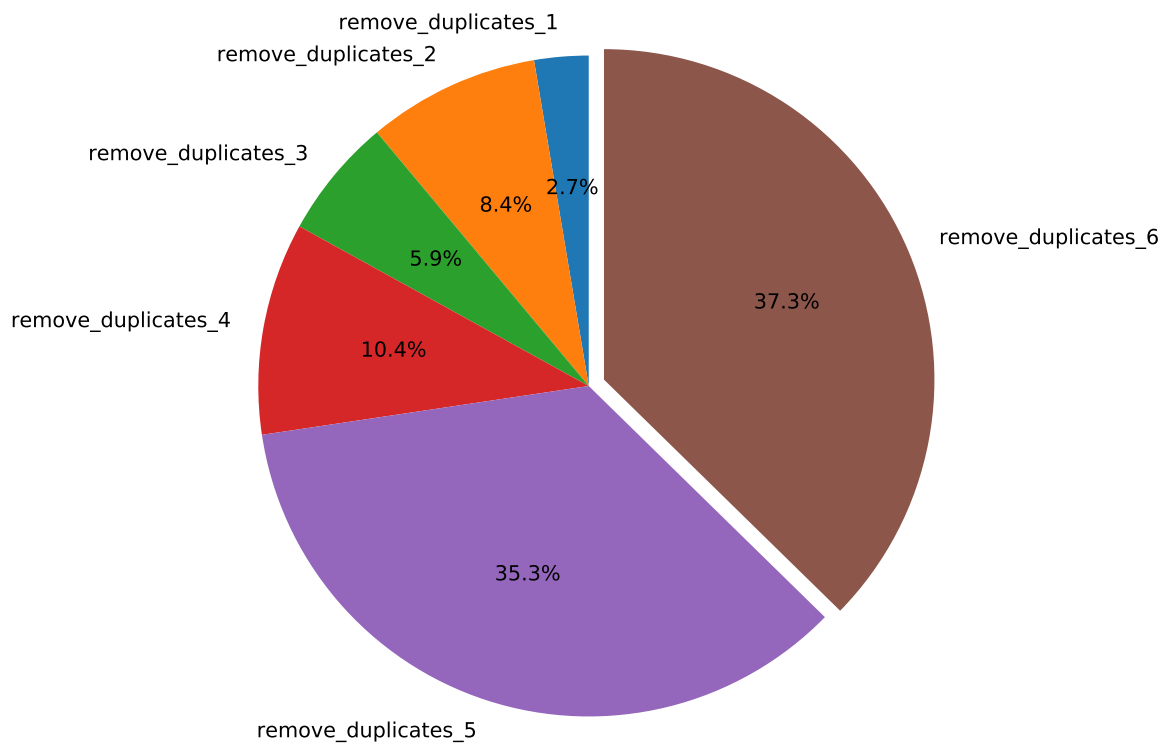
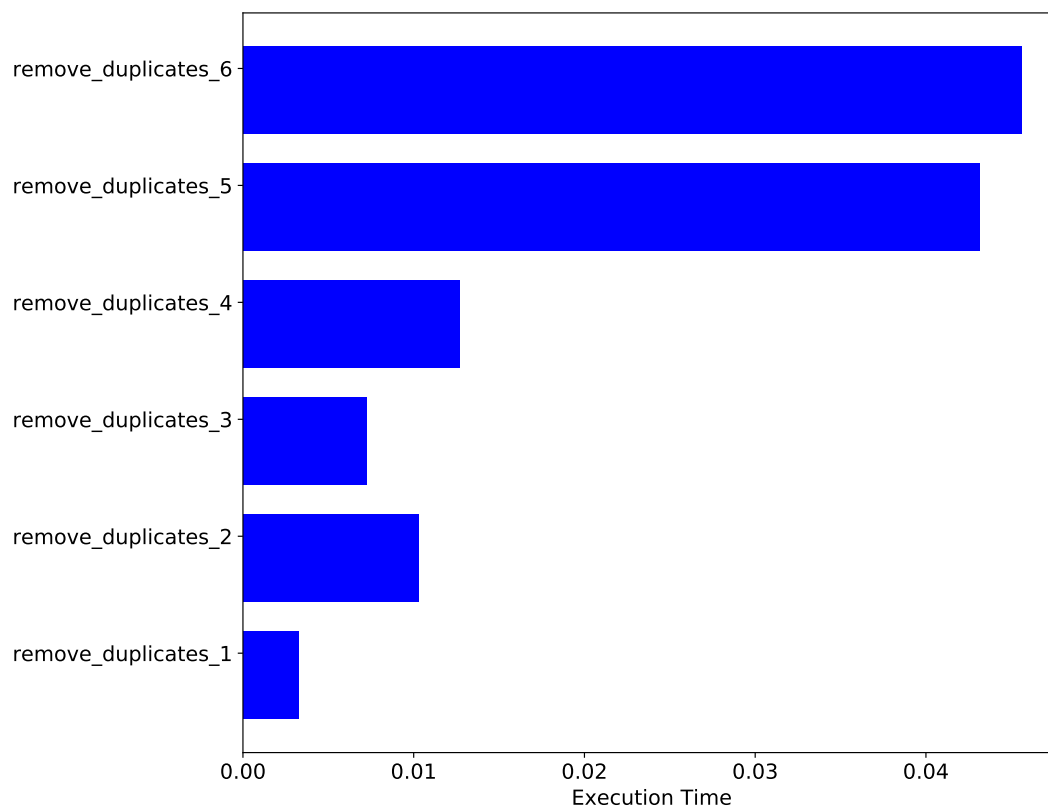
TIMES FOR THE FIRST 78000 NUMBERS



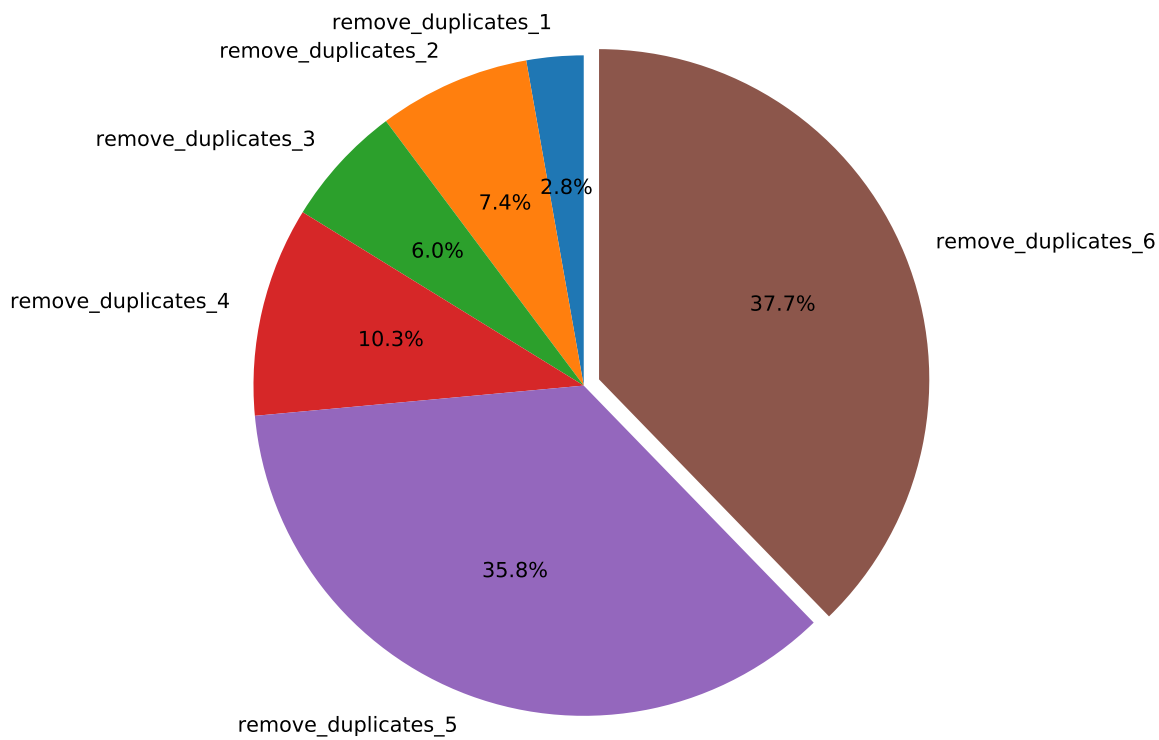
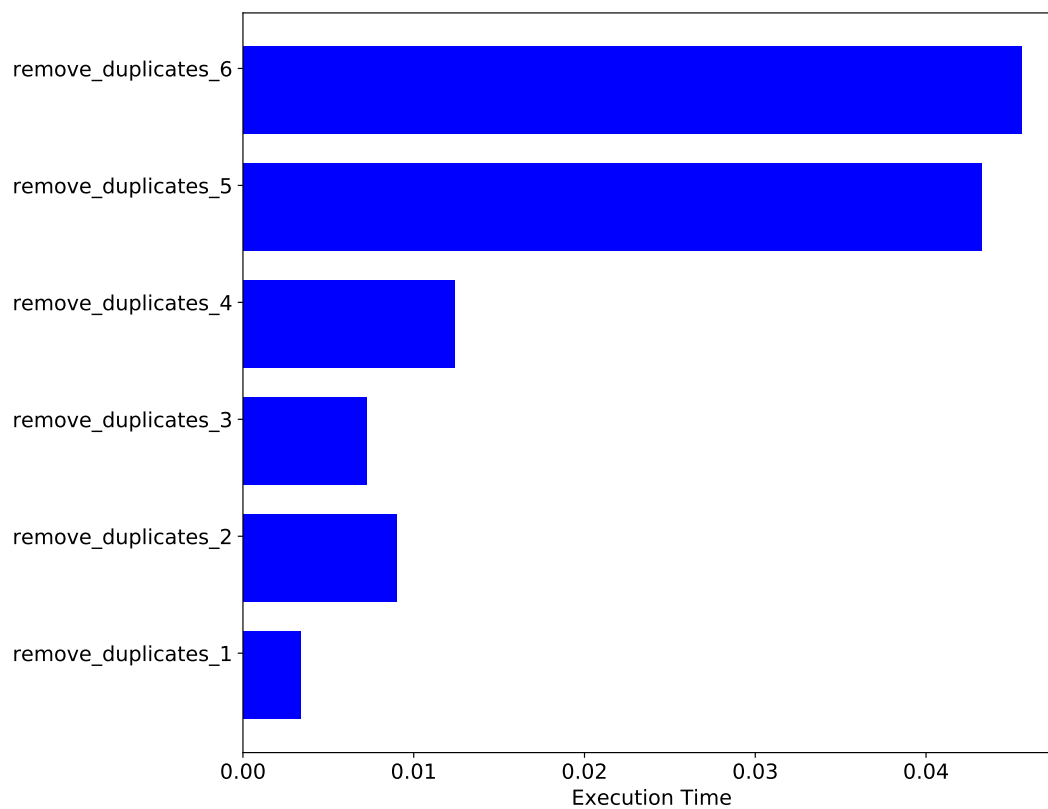
TIMES FOR THE FIRST 80000 NUMBERS



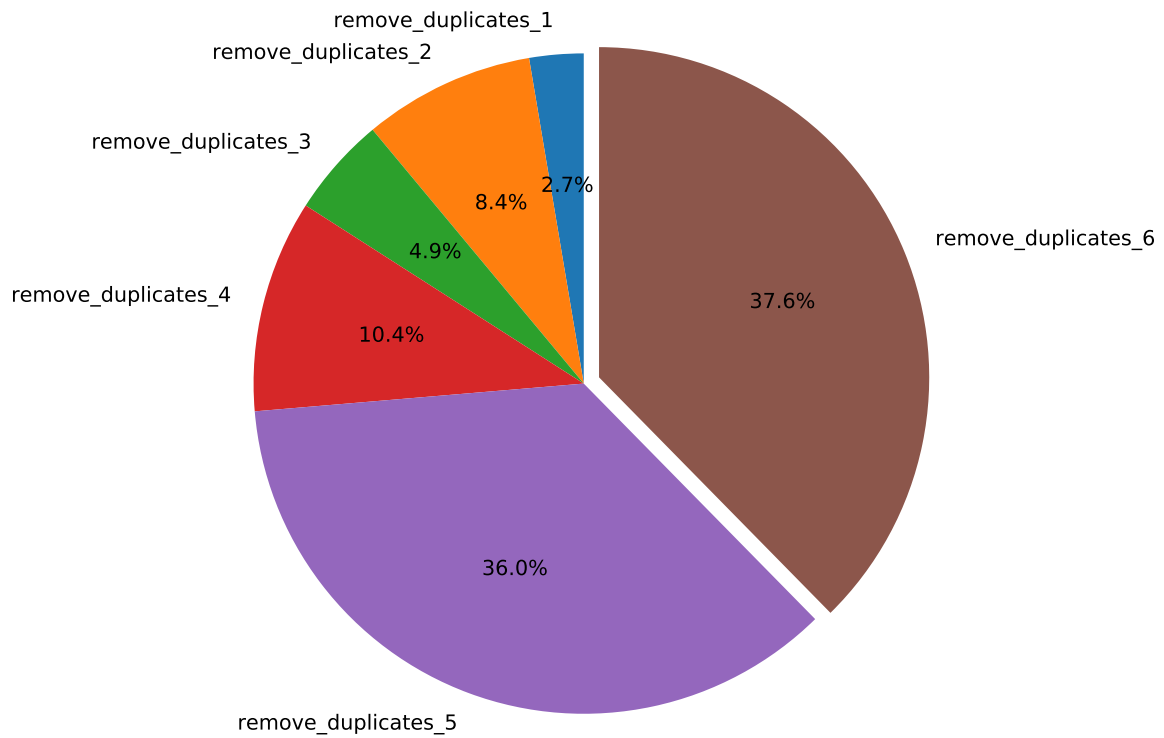
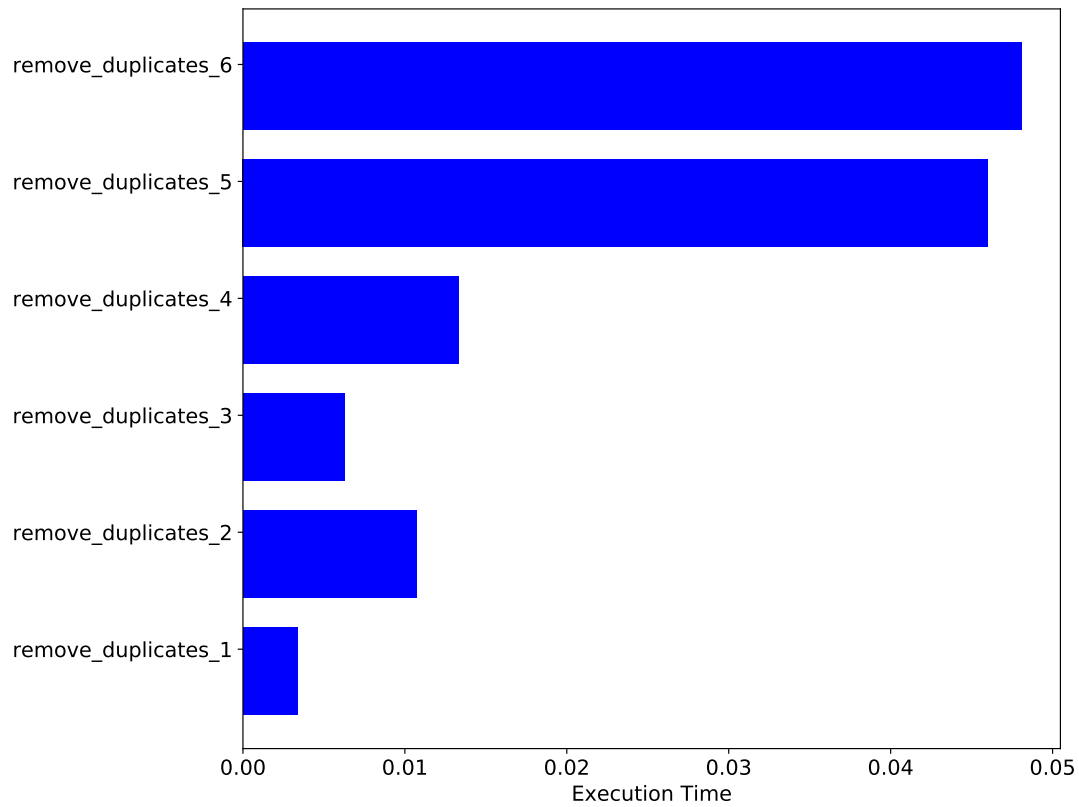
TIMES FOR THE FIRST 82000 NUMBERS



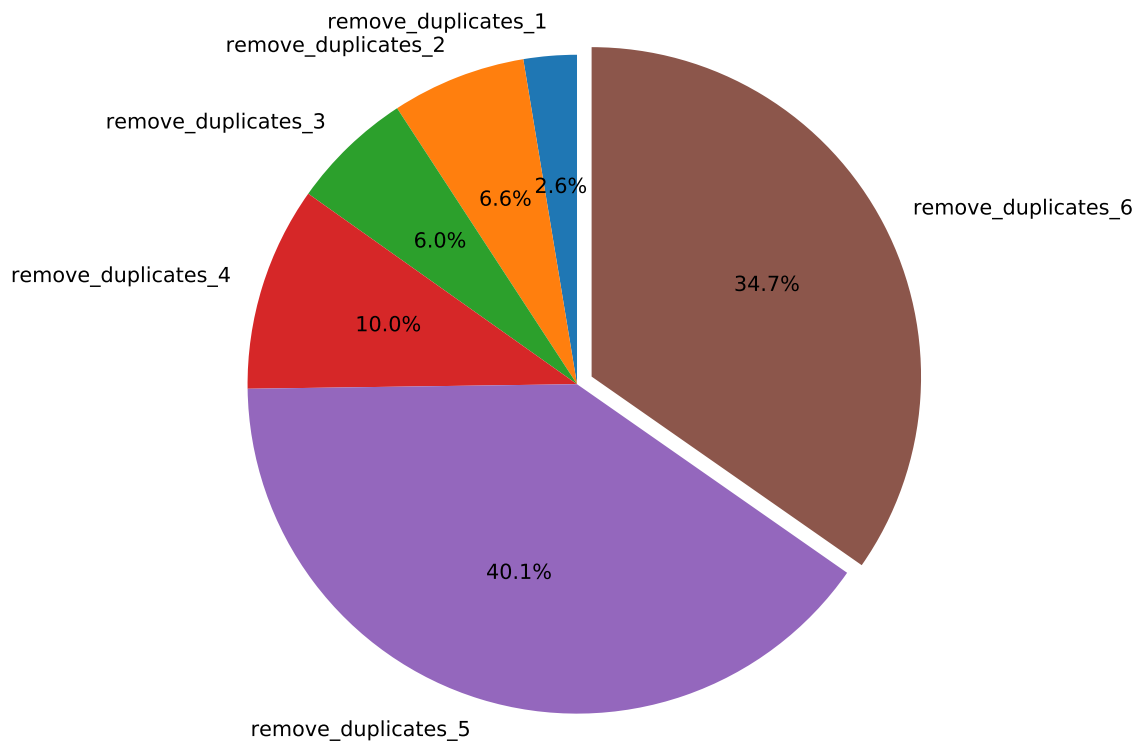
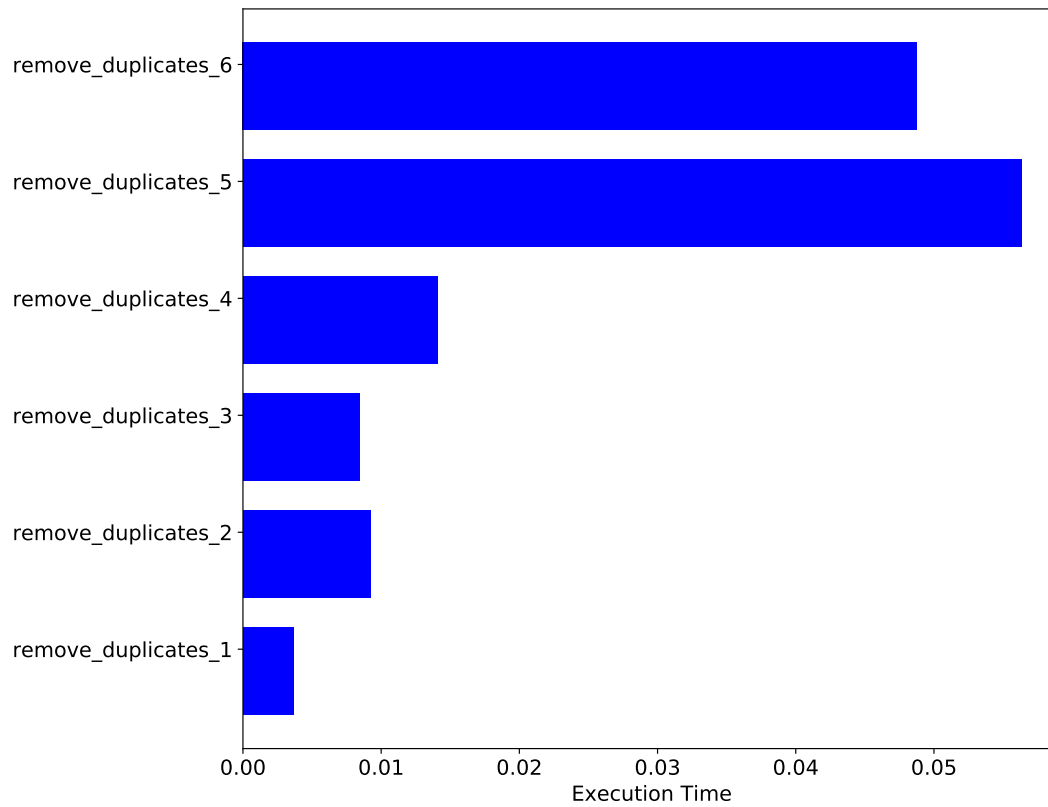
TIMES FOR THE FIRST 84000 NUMBERS



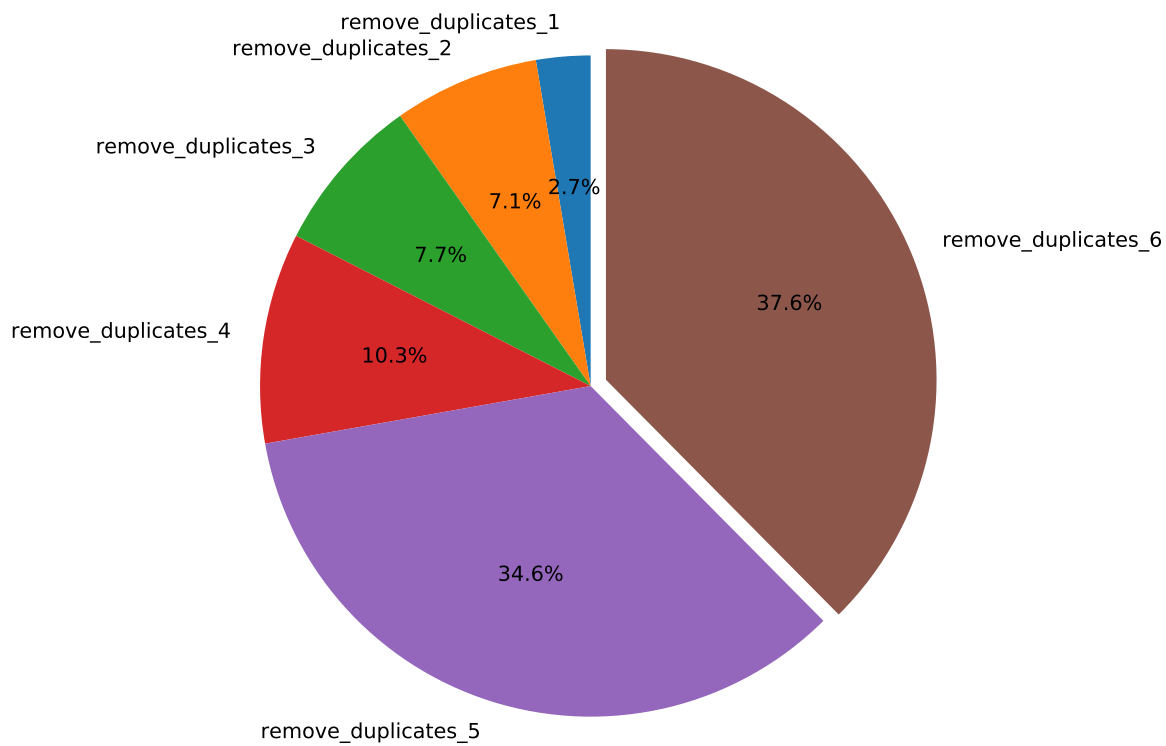
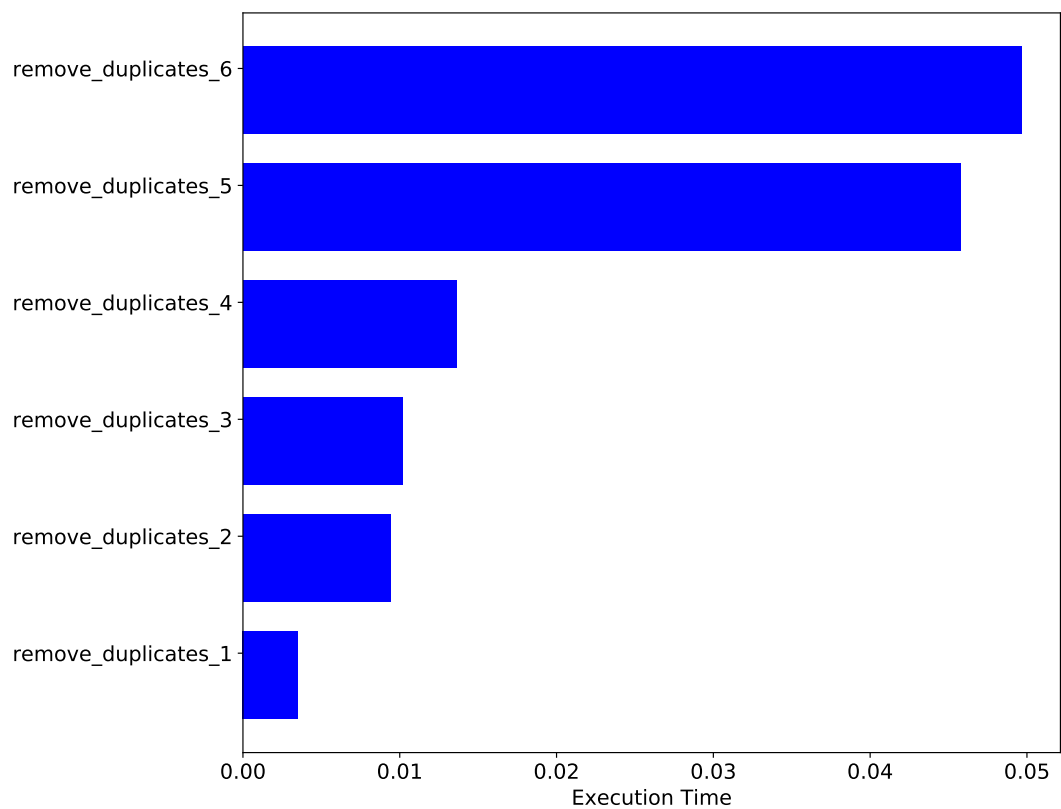
TIMES FOR THE FIRST 86000 NUMBERS



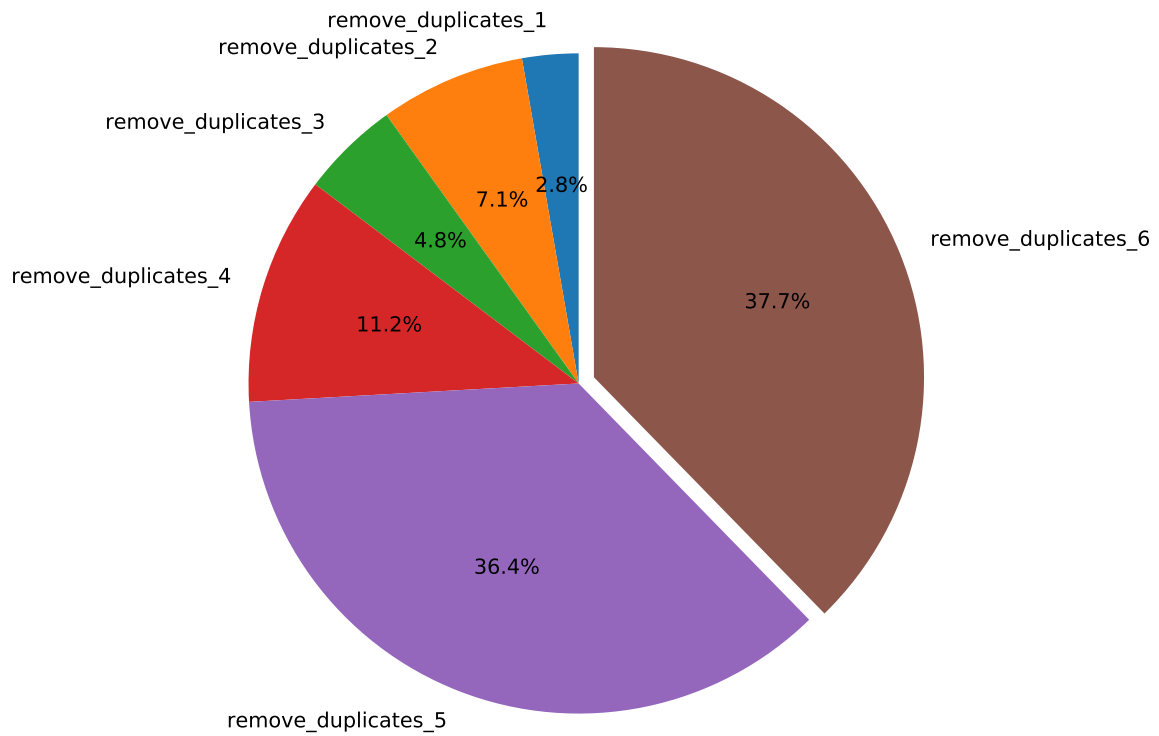
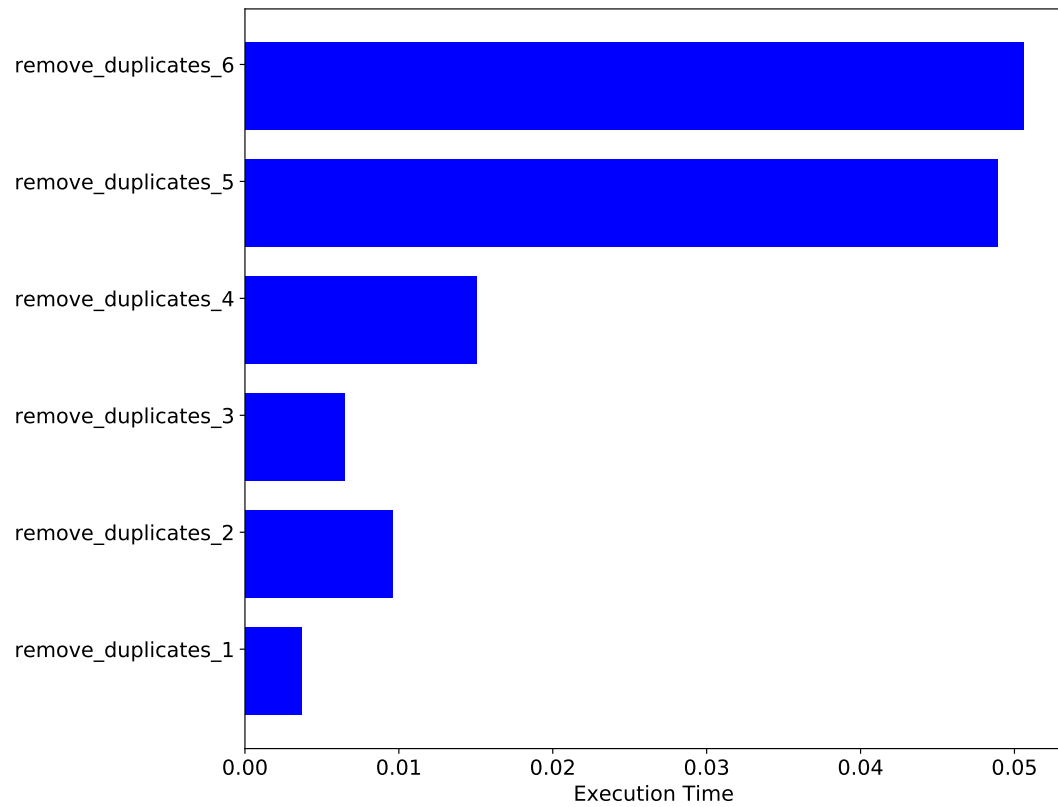
TIMES FOR THE FIRST 88000 NUMBERS



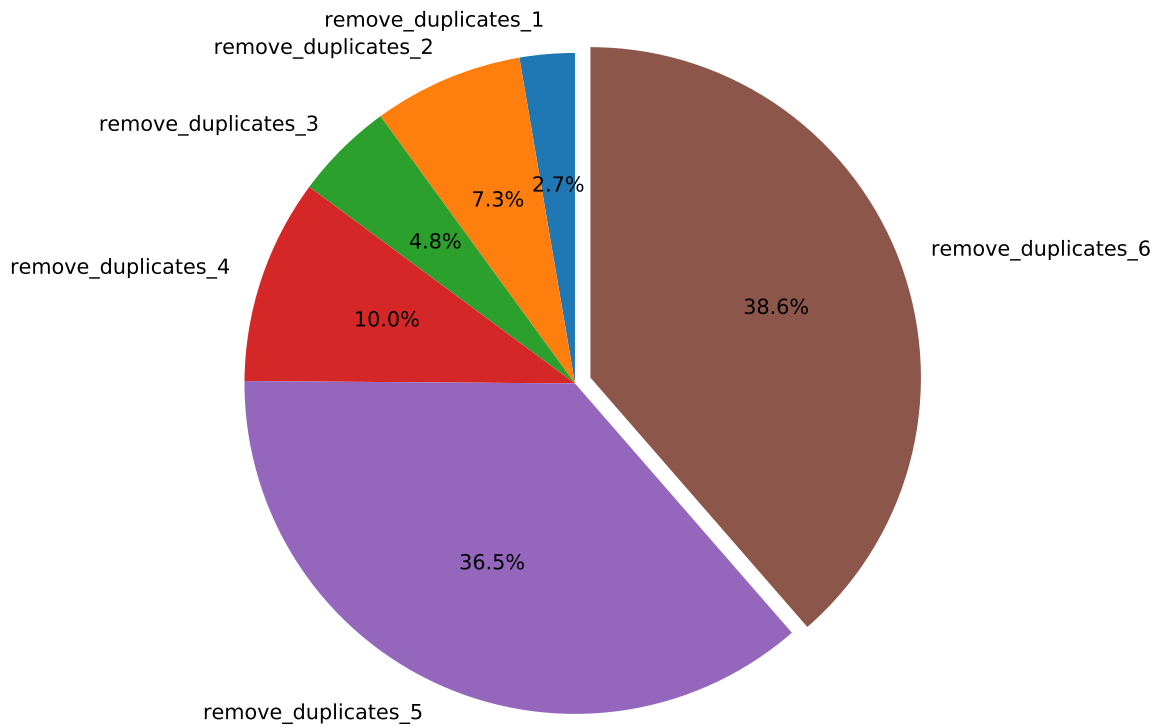
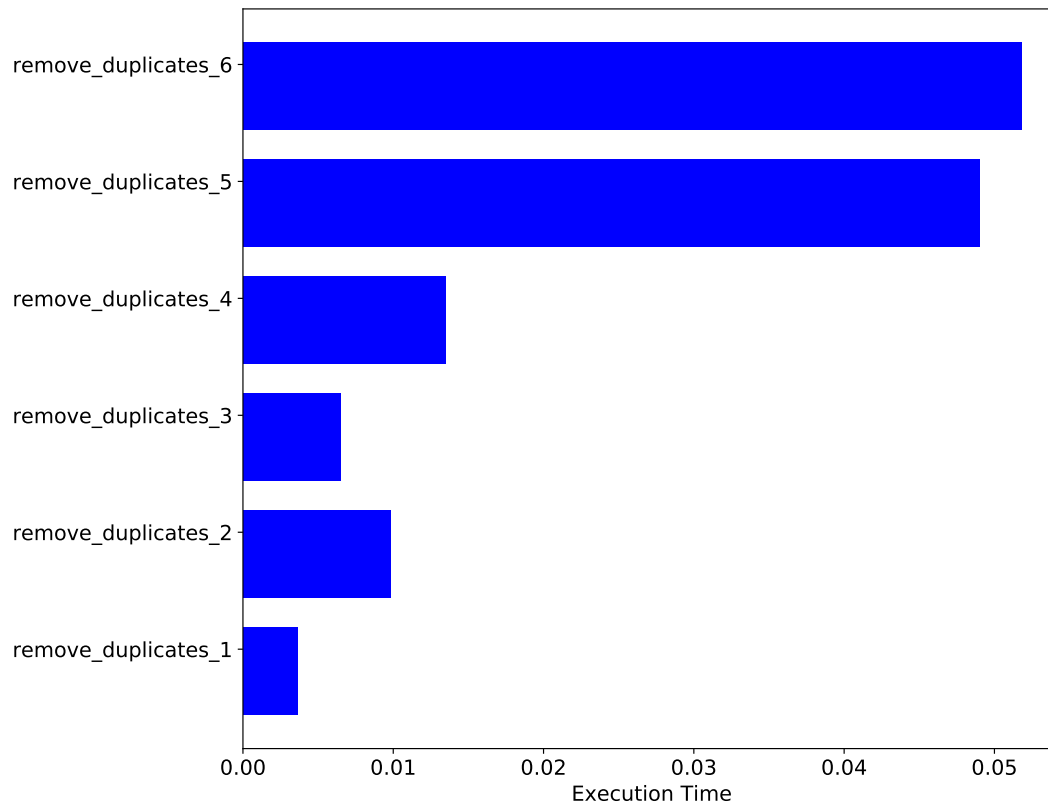
TIMES FOR THE FIRST 90000 NUMBERS



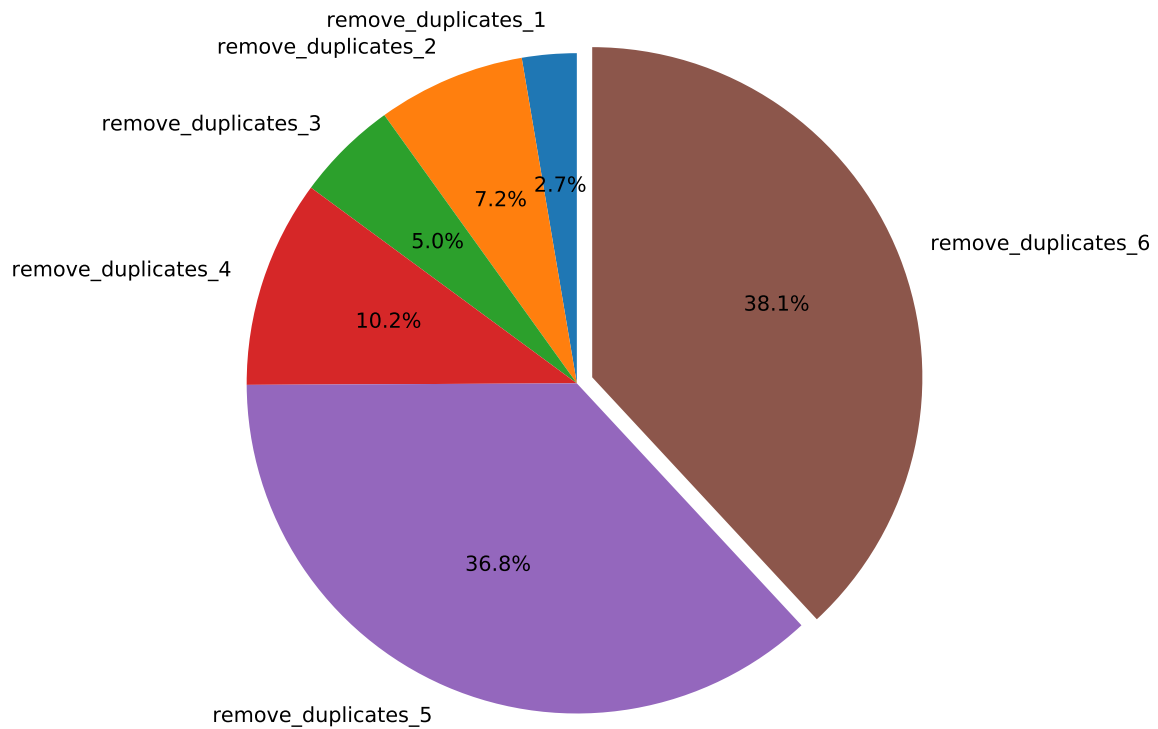
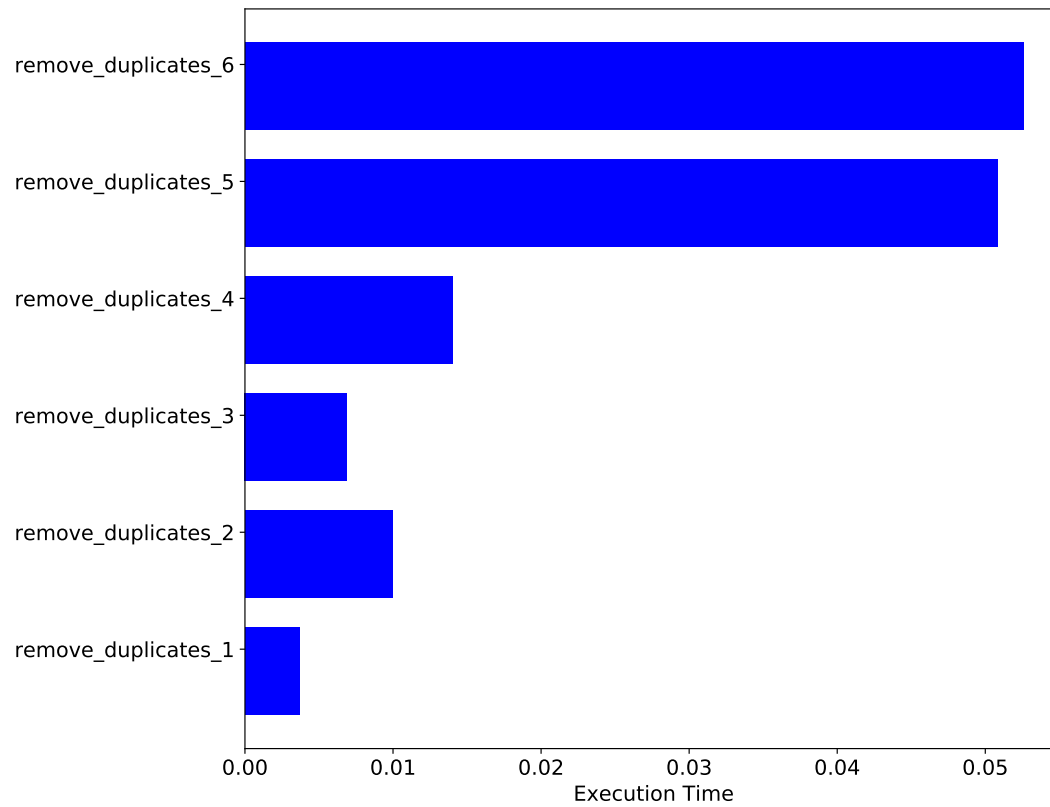
TIMES FOR THE FIRST 92000 NUMBERS



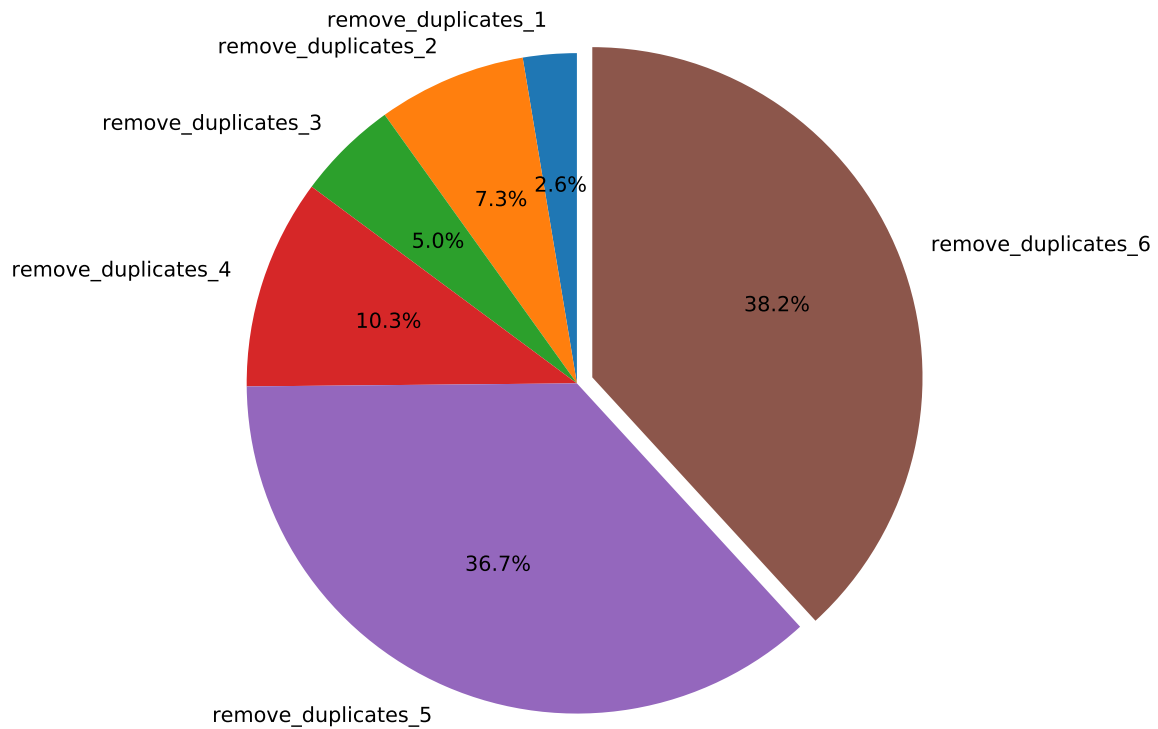
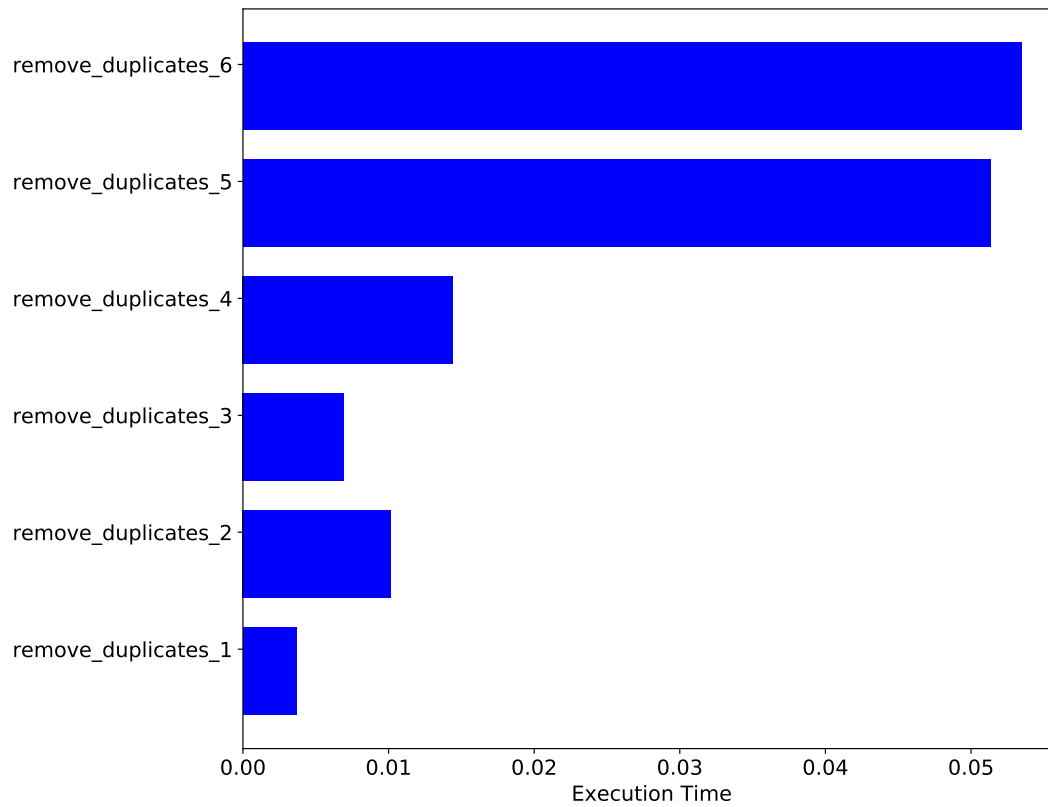
TIMES FOR THE FIRST 94000 NUMBERS



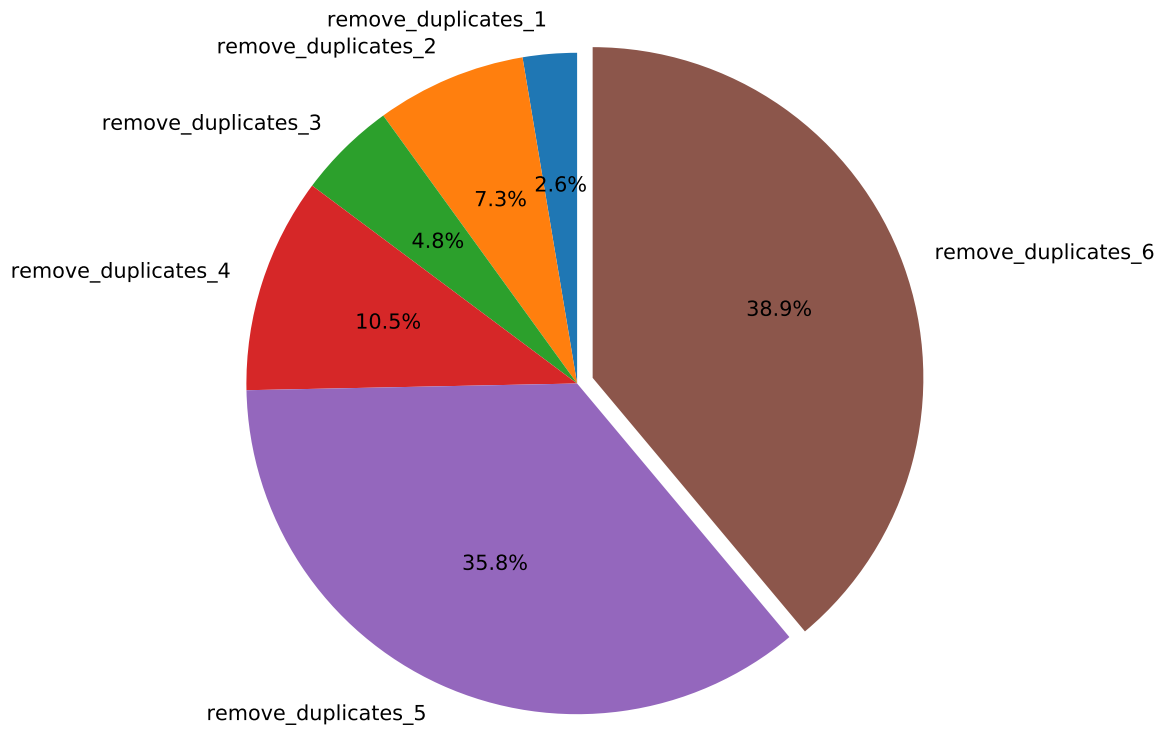
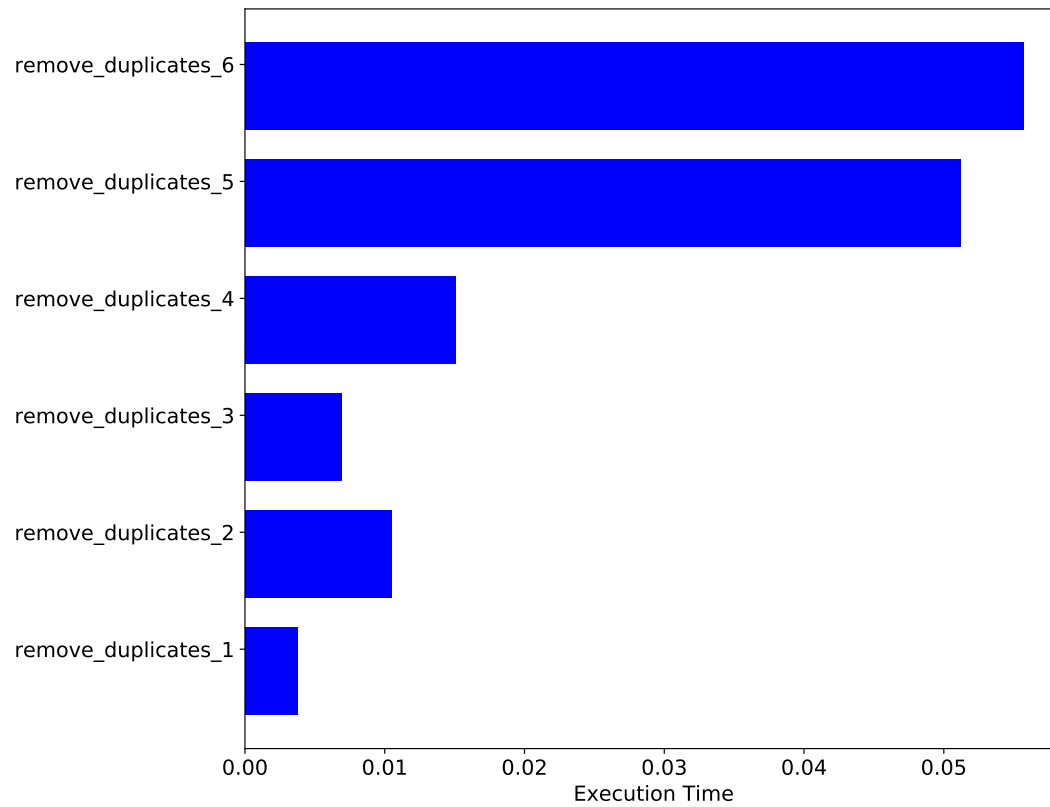
TIMES FOR THE FIRST 96000 NUMBERS



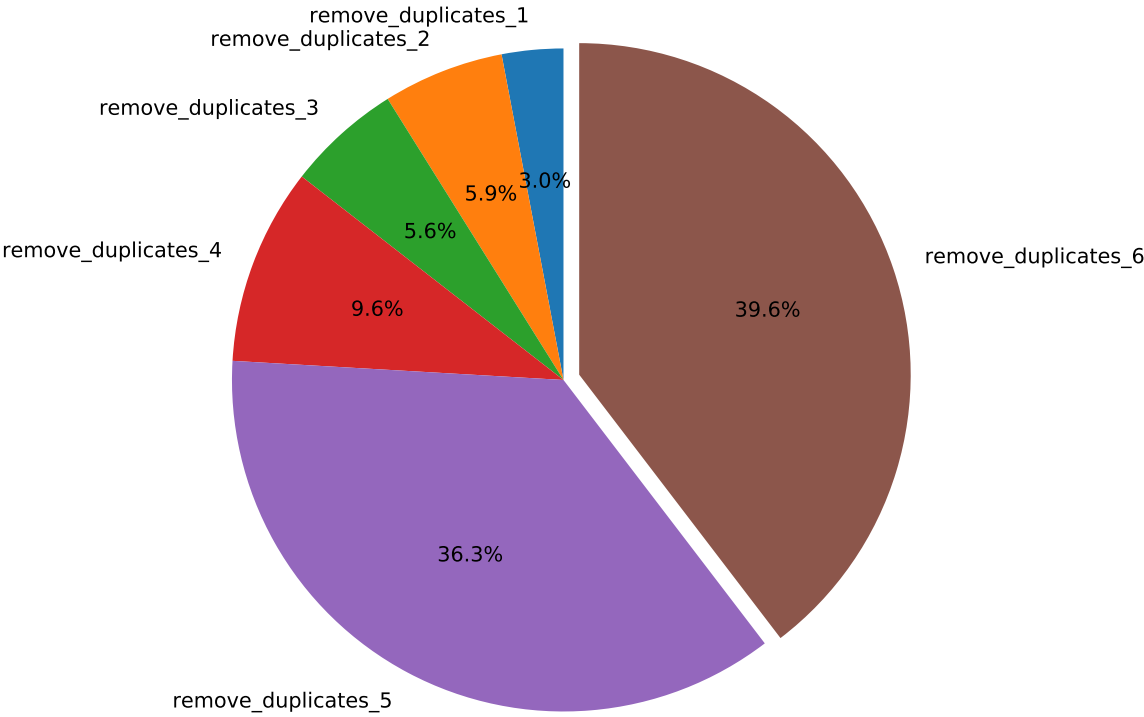
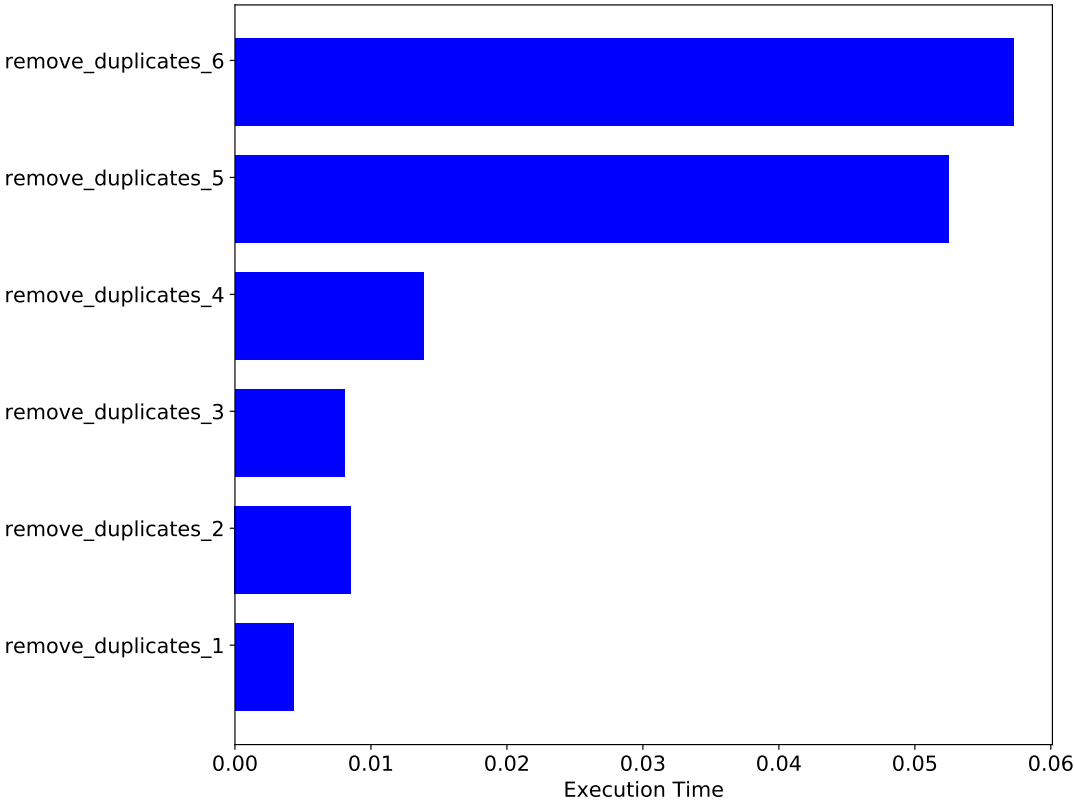
TIMES FOR THE FIRST 98000 NUMBERS



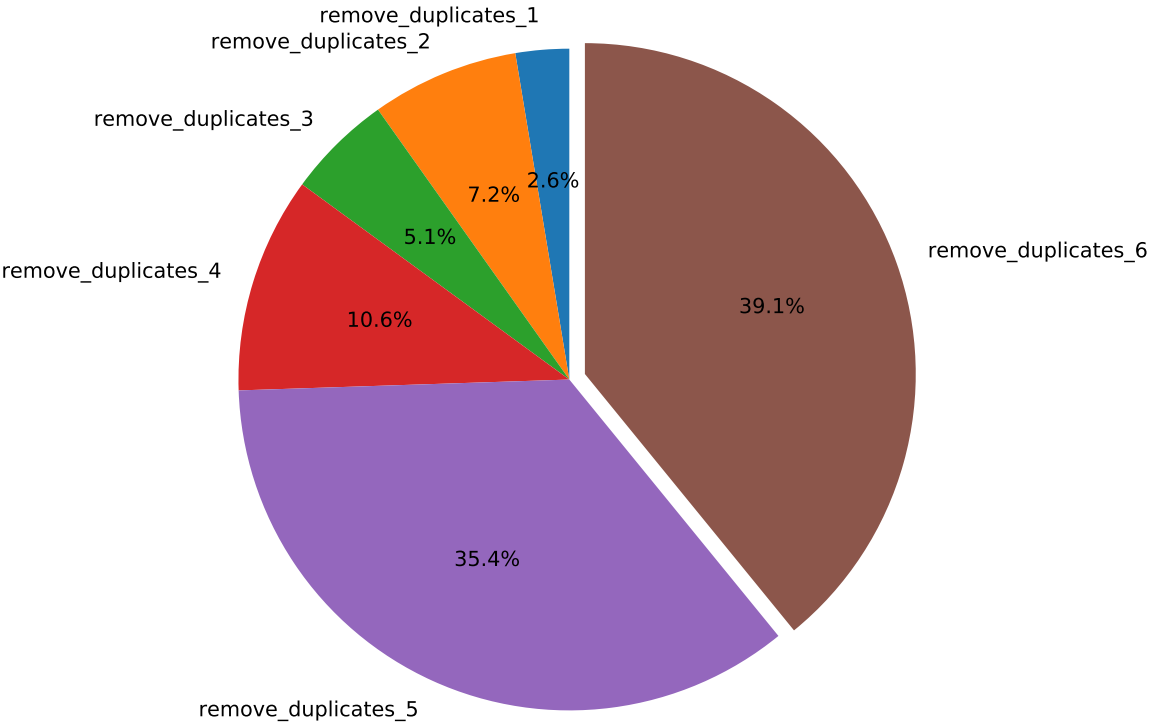
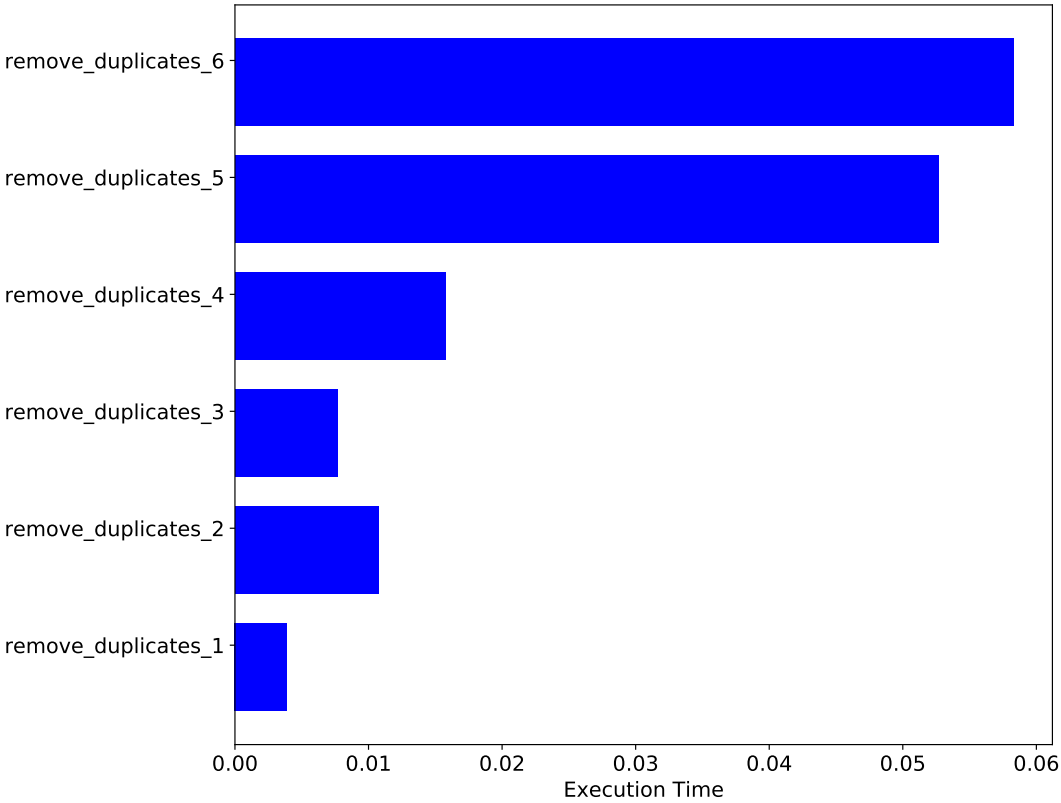
TIMES FOR THE FIRST 100000 NUMBERS



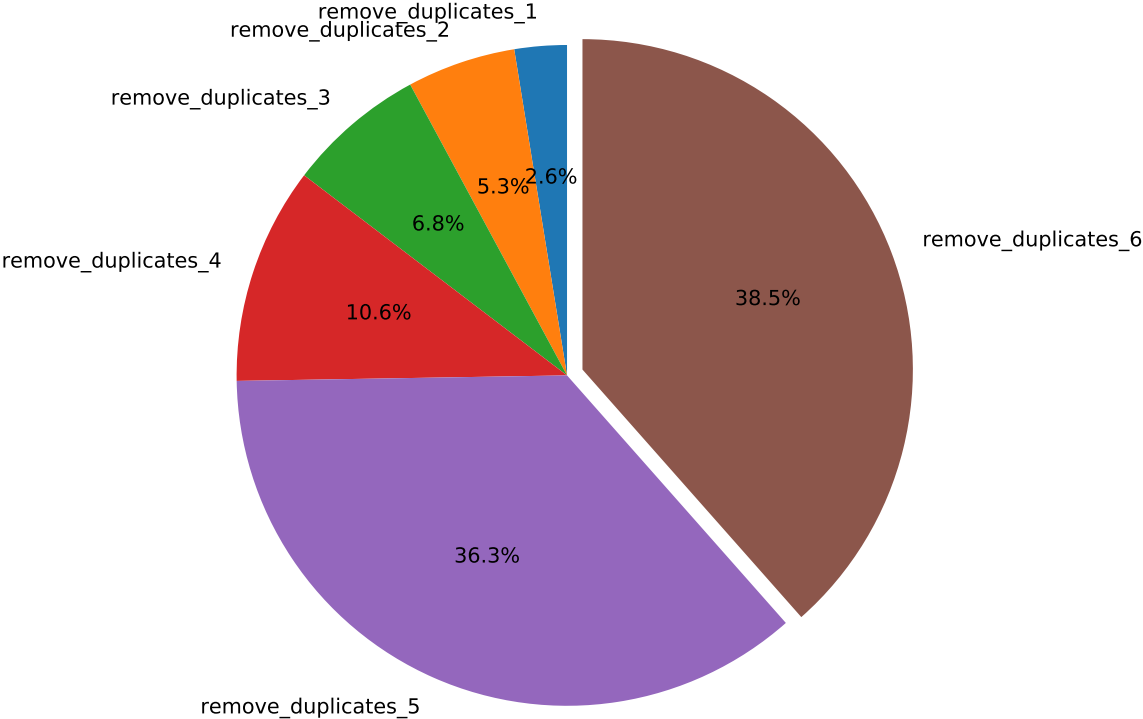
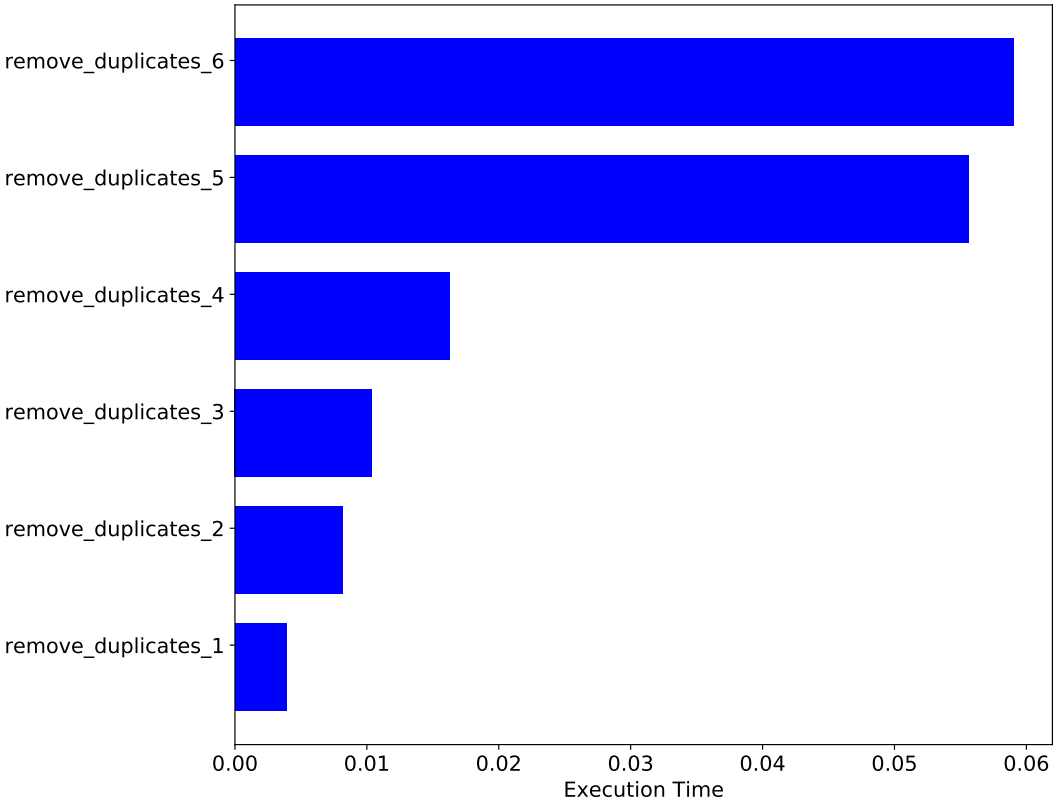
TIMES FOR THE FIRST 102000 NUMBERS



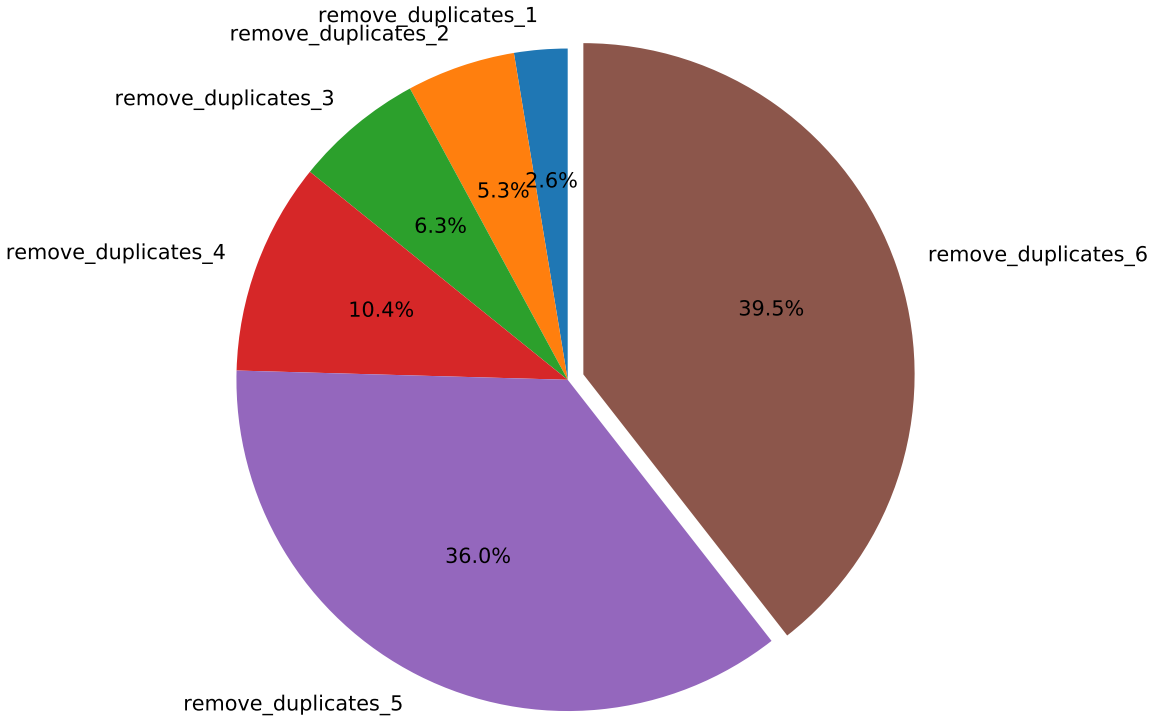
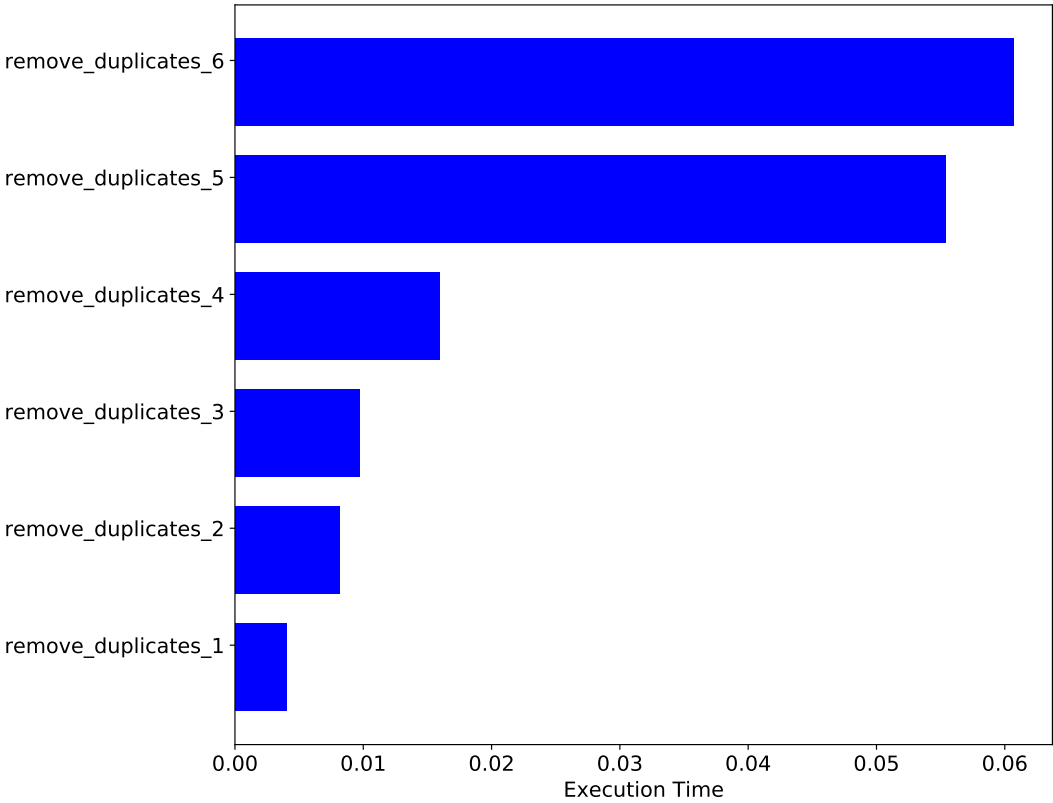
TIMES FOR THE FIRST 104000 NUMBERS



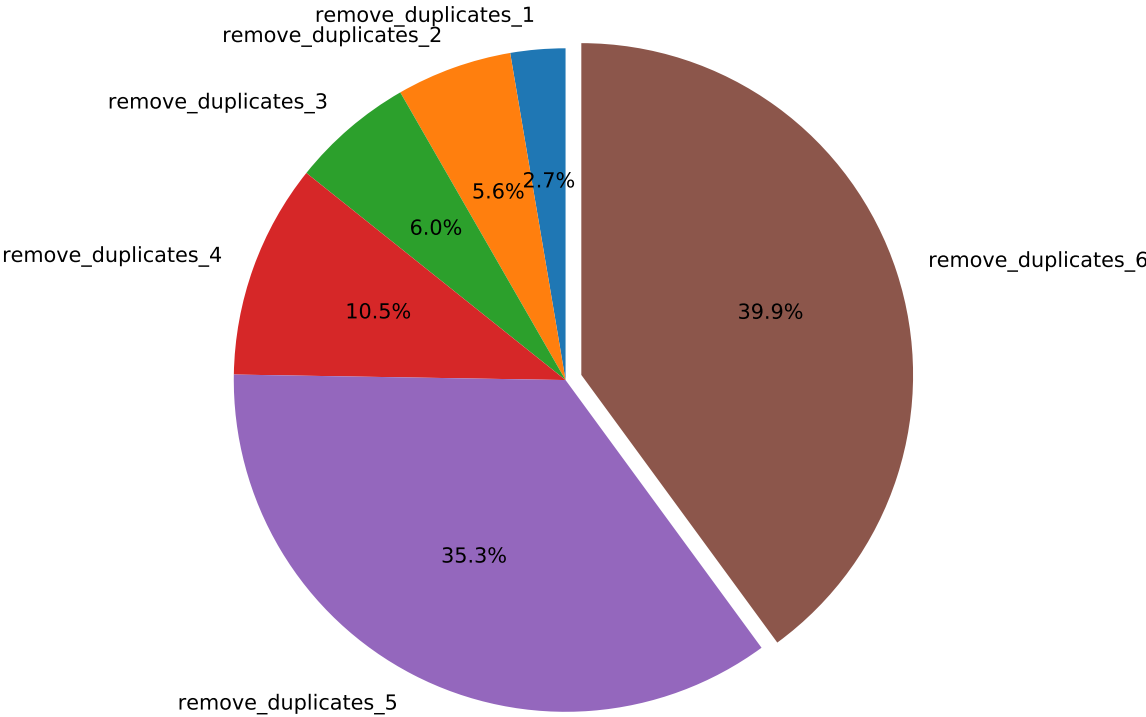
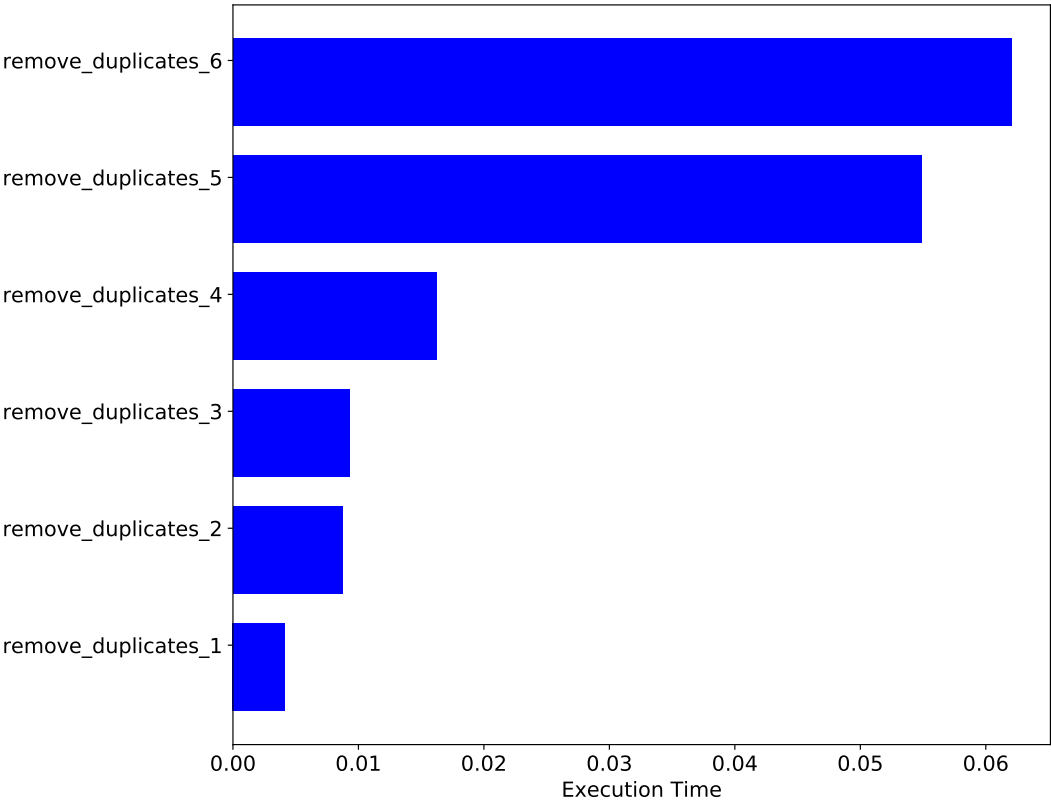
TIMES FOR THE FIRST 106000 NUMBERS



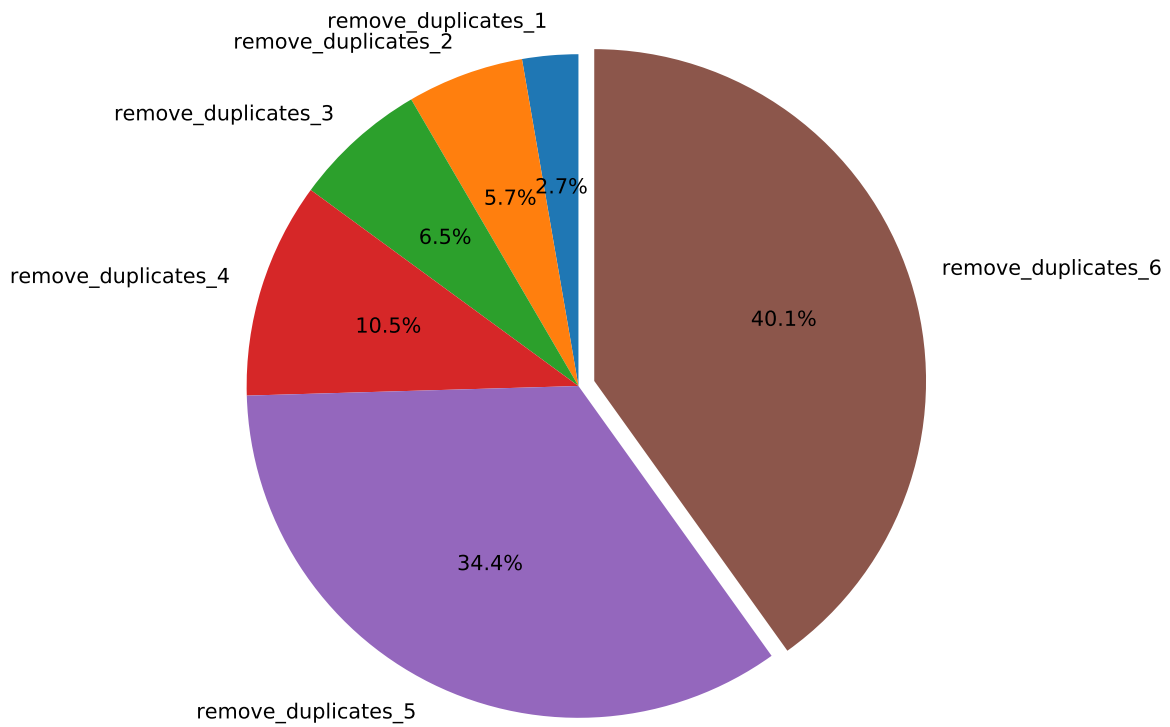
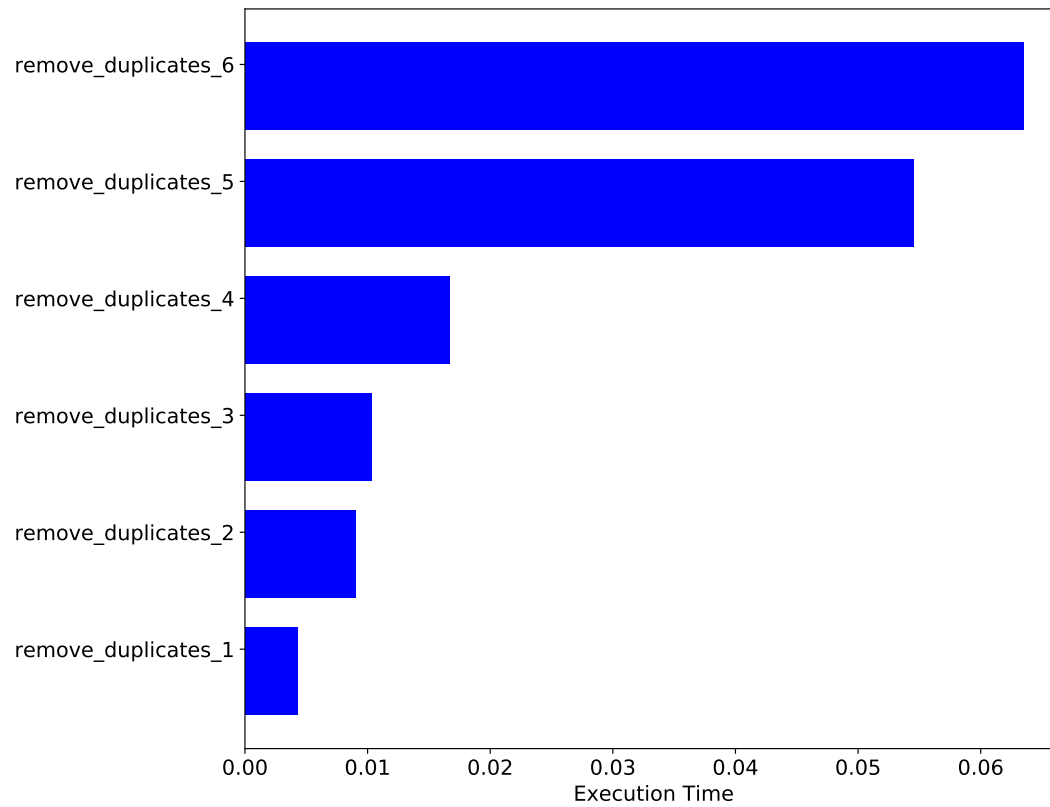
TIMES FOR THE FIRST 108000 NUMBERS



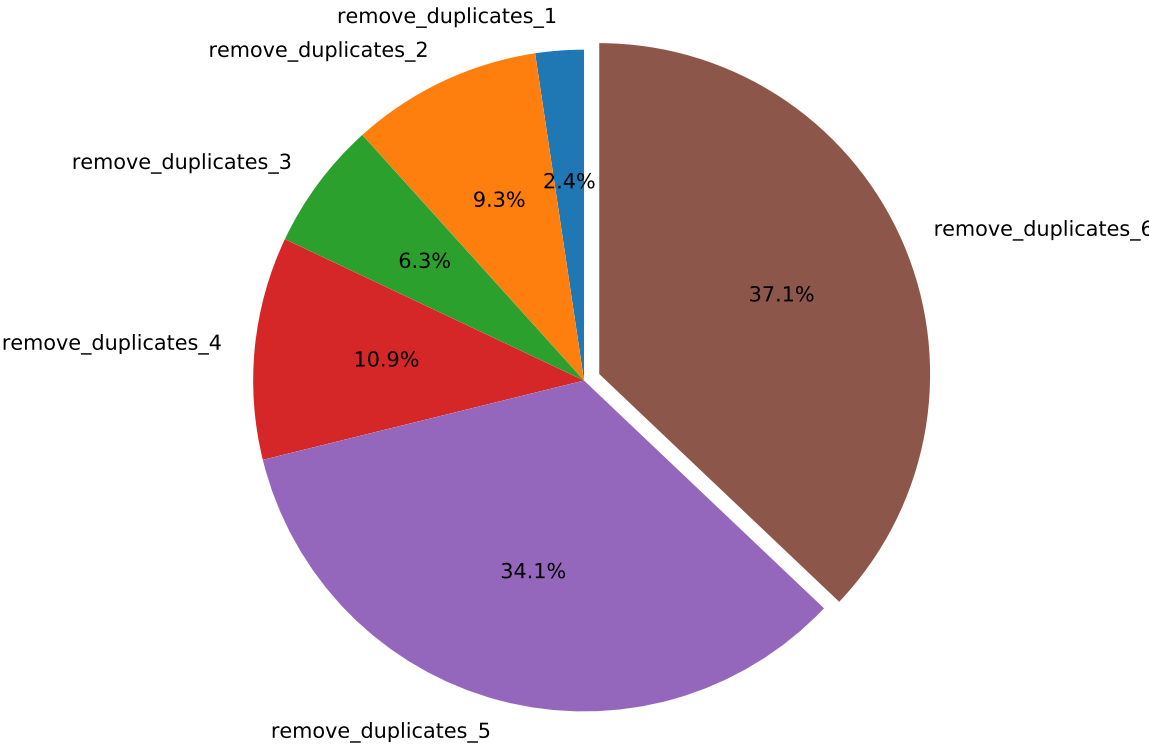
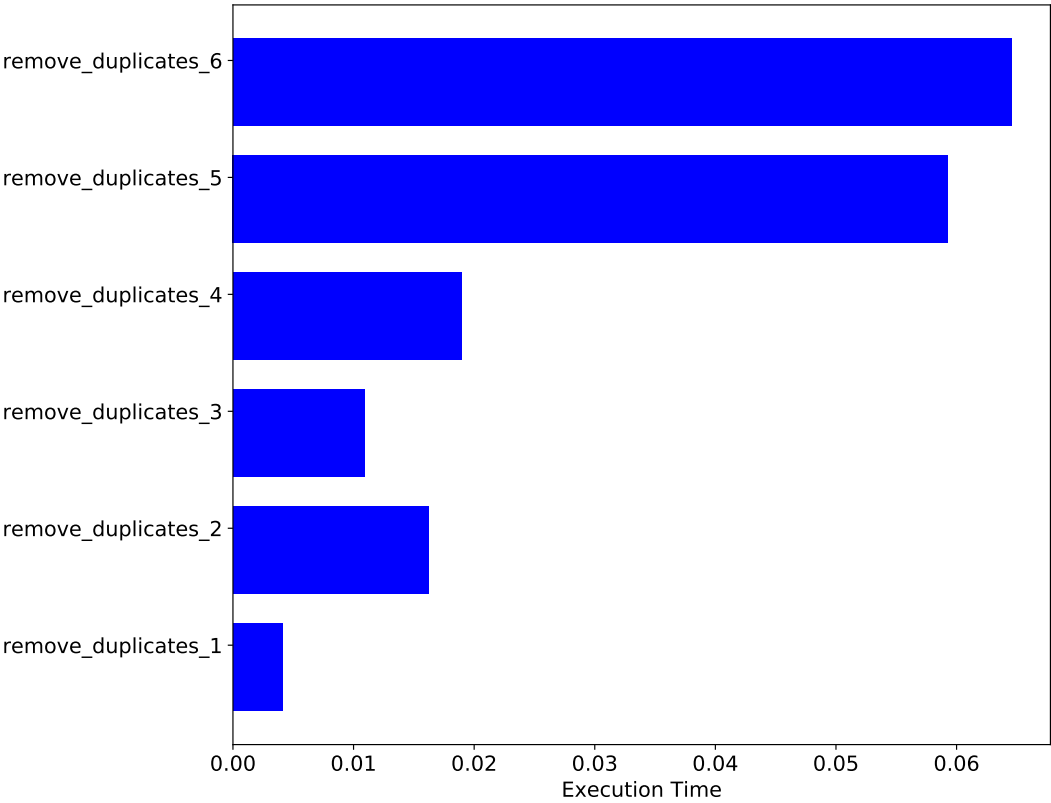
TIMES FOR THE FIRST 110000 NUMBERS



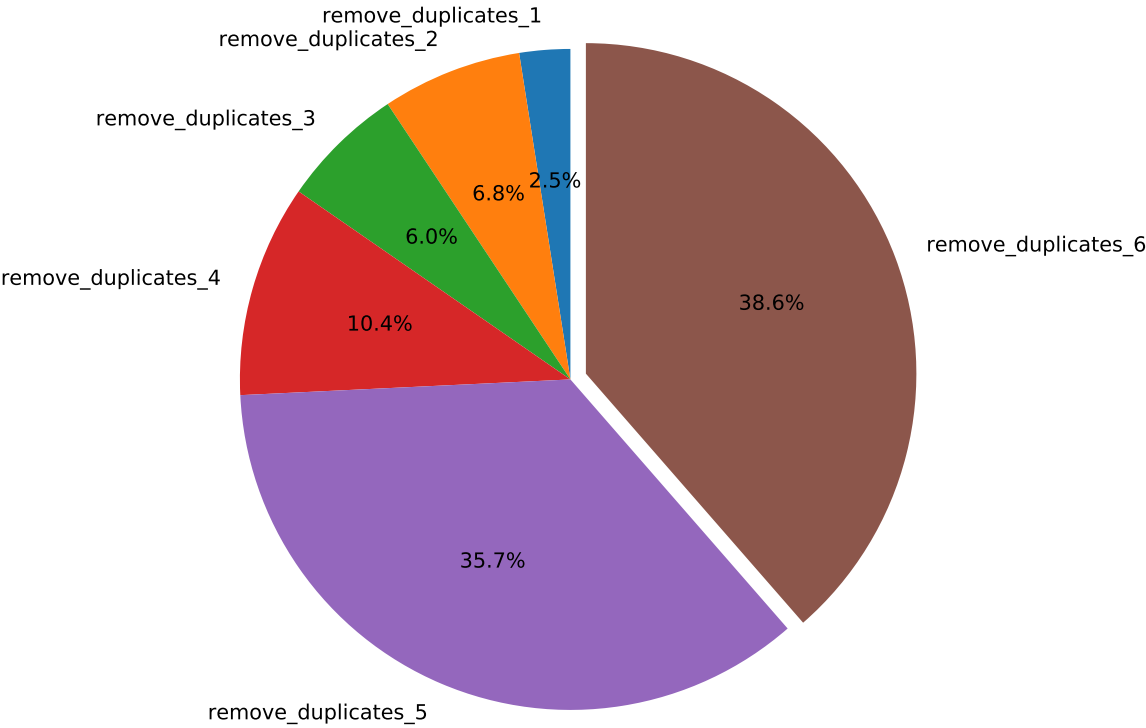
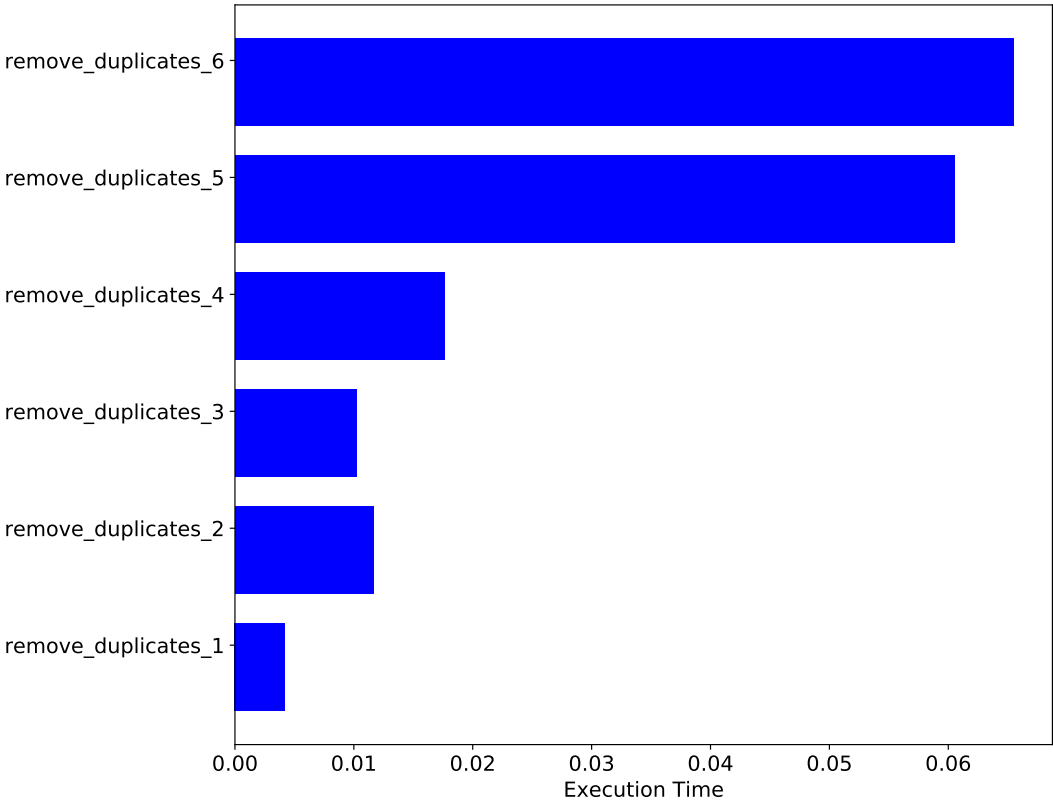
TIMES FOR THE FIRST 112000 NUMBERS



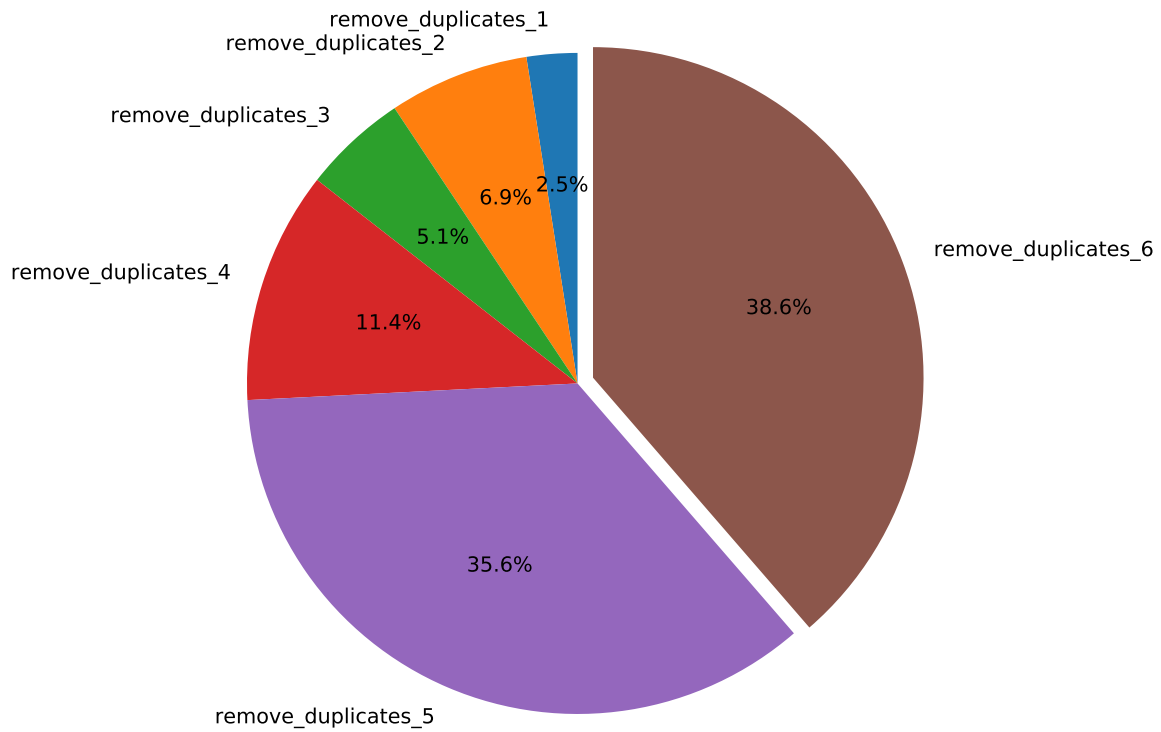
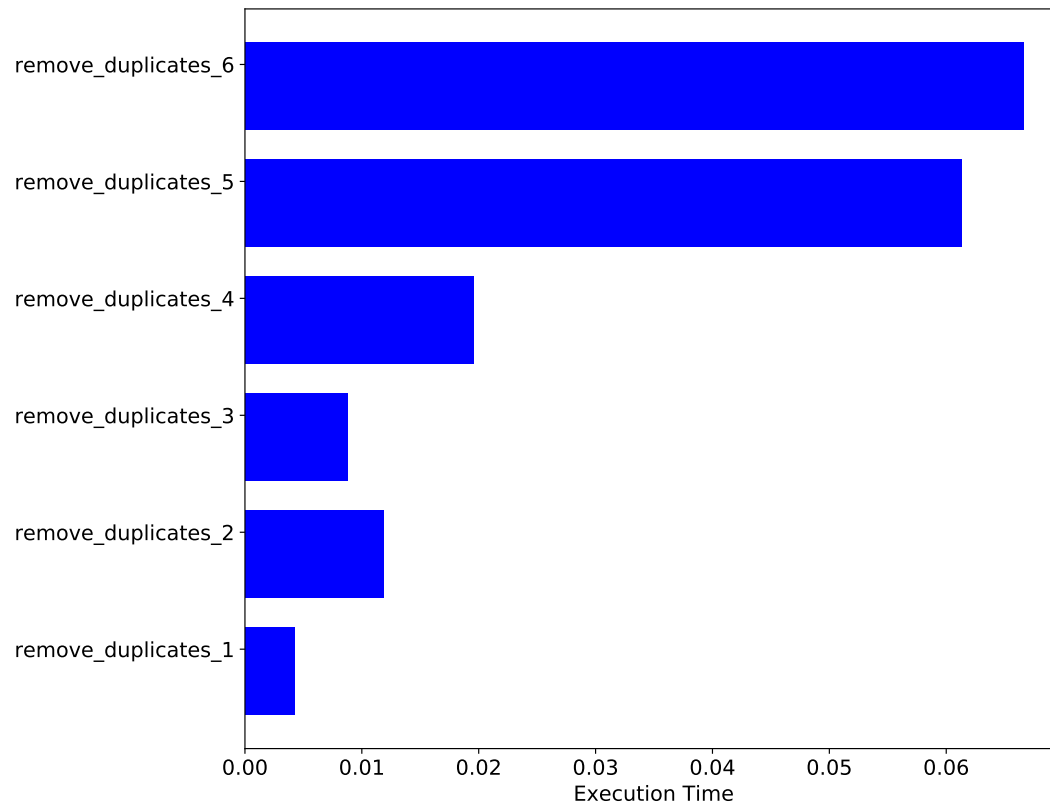
TIMES FOR THE FIRST 114000 NUMBERS



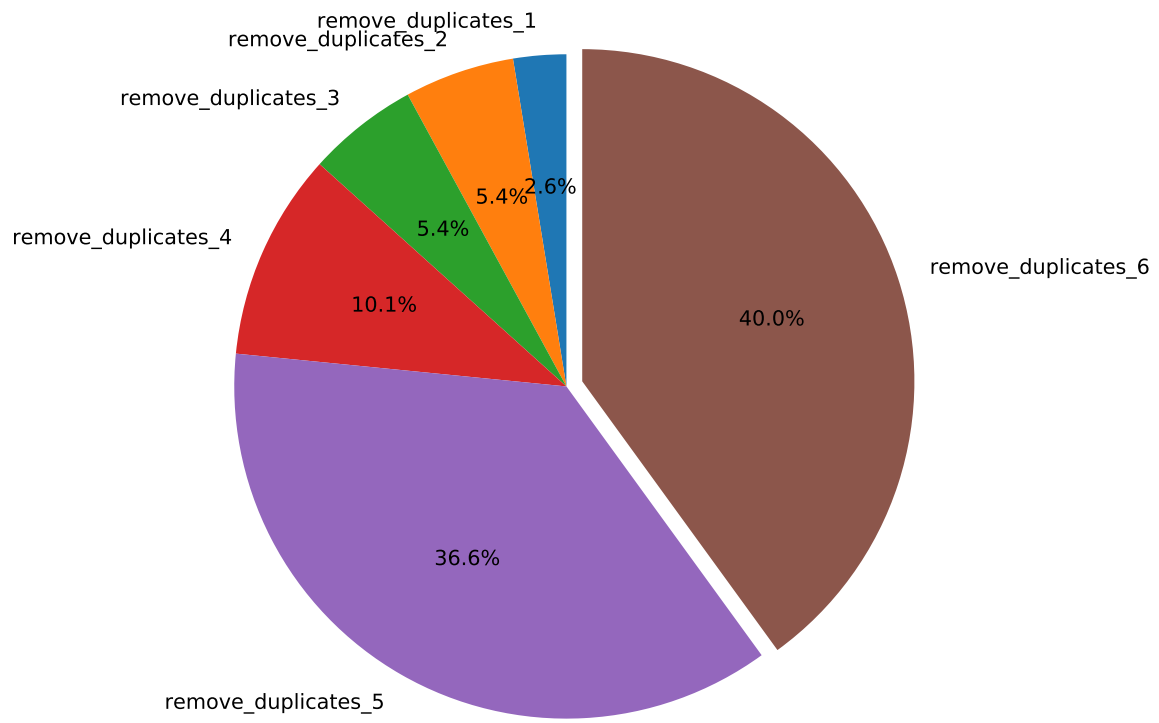
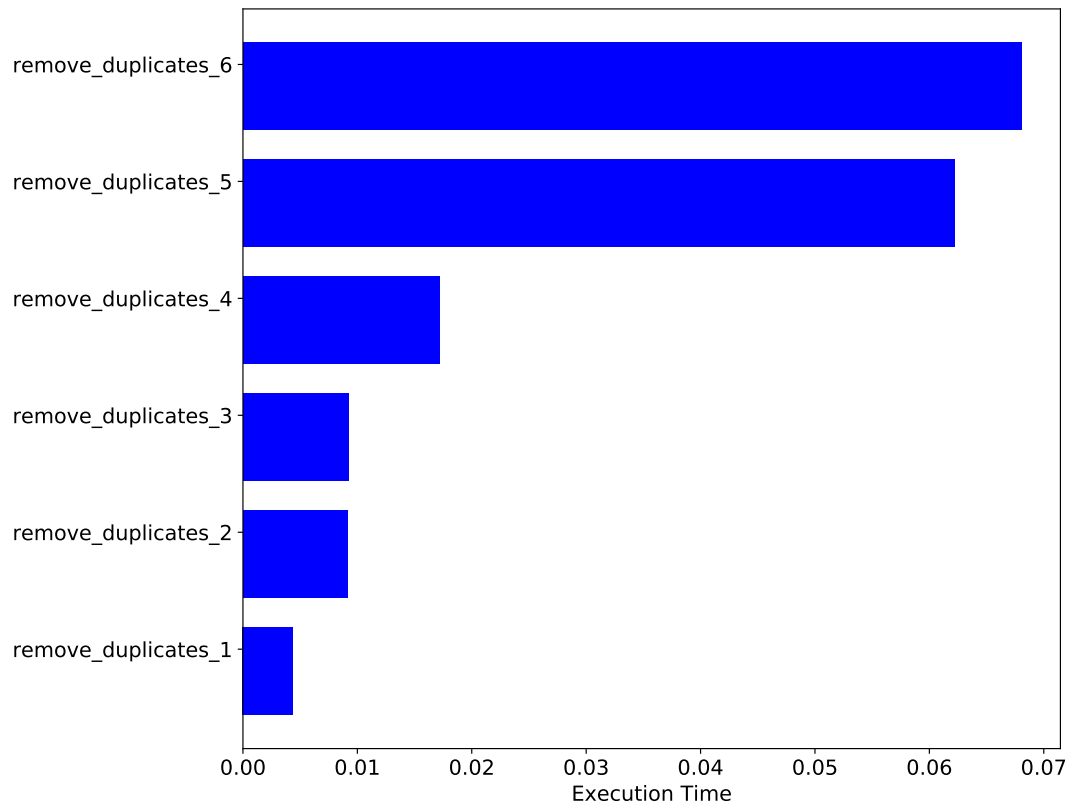
TIMES FOR THE FIRST 116000 NUMBERS



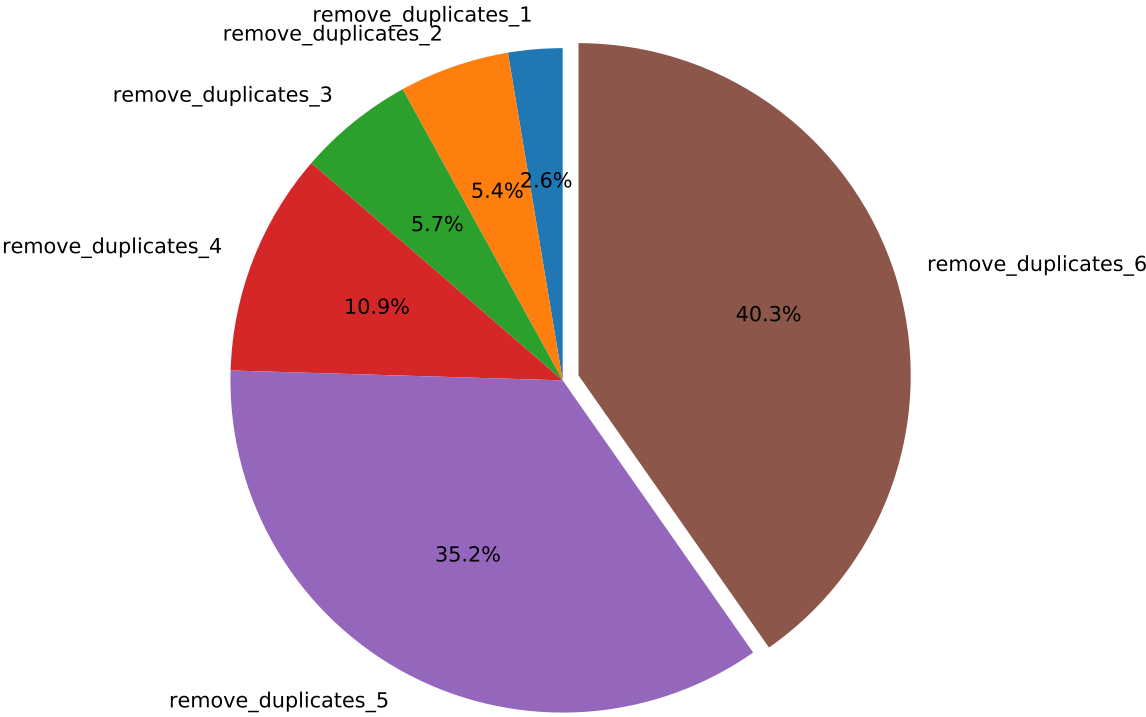
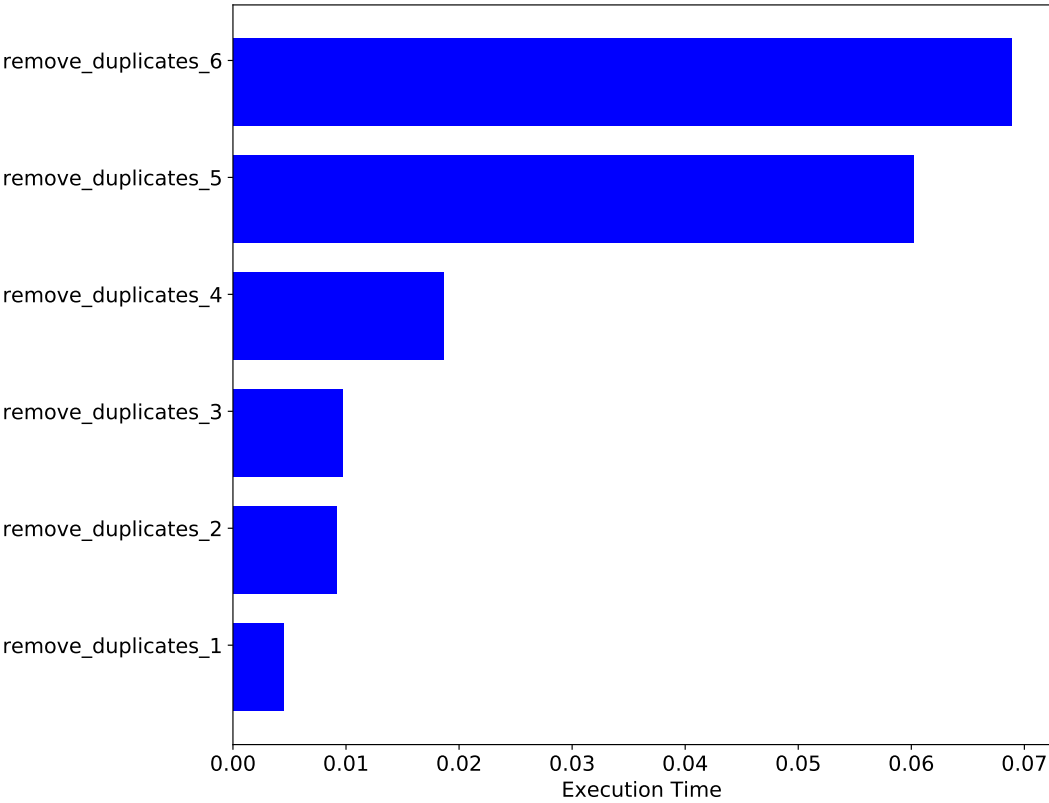
TIMES FOR THE FIRST 118000 NUMBERS



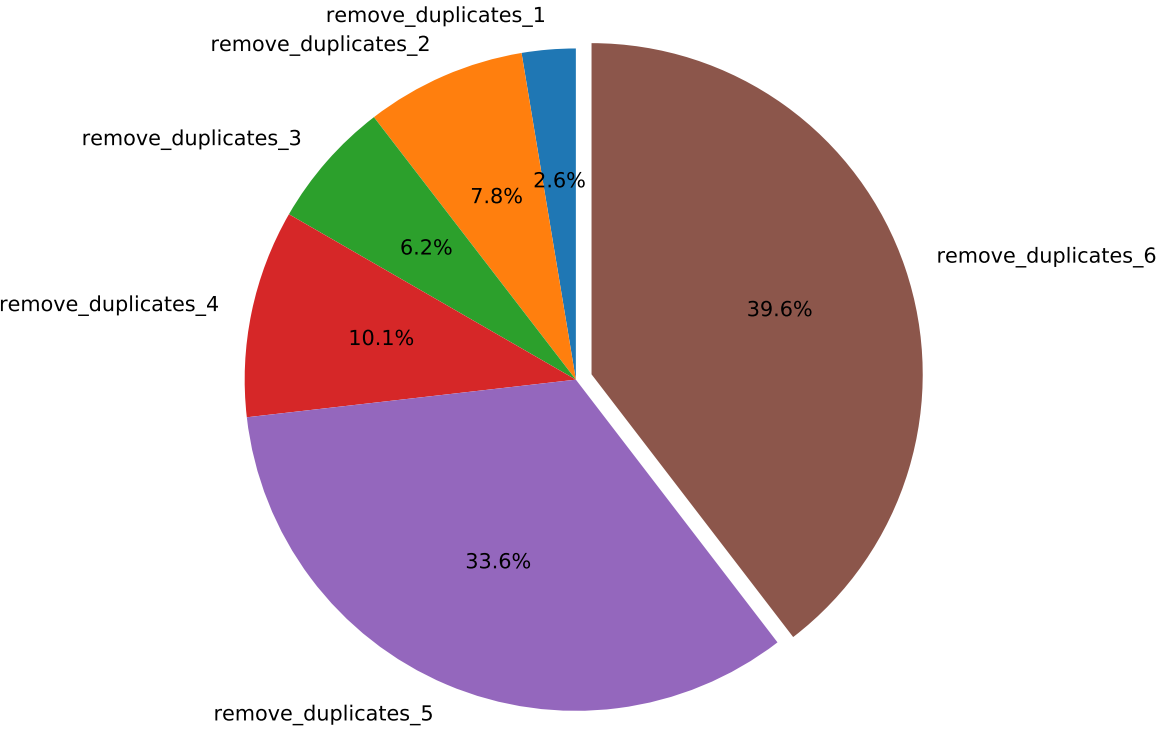
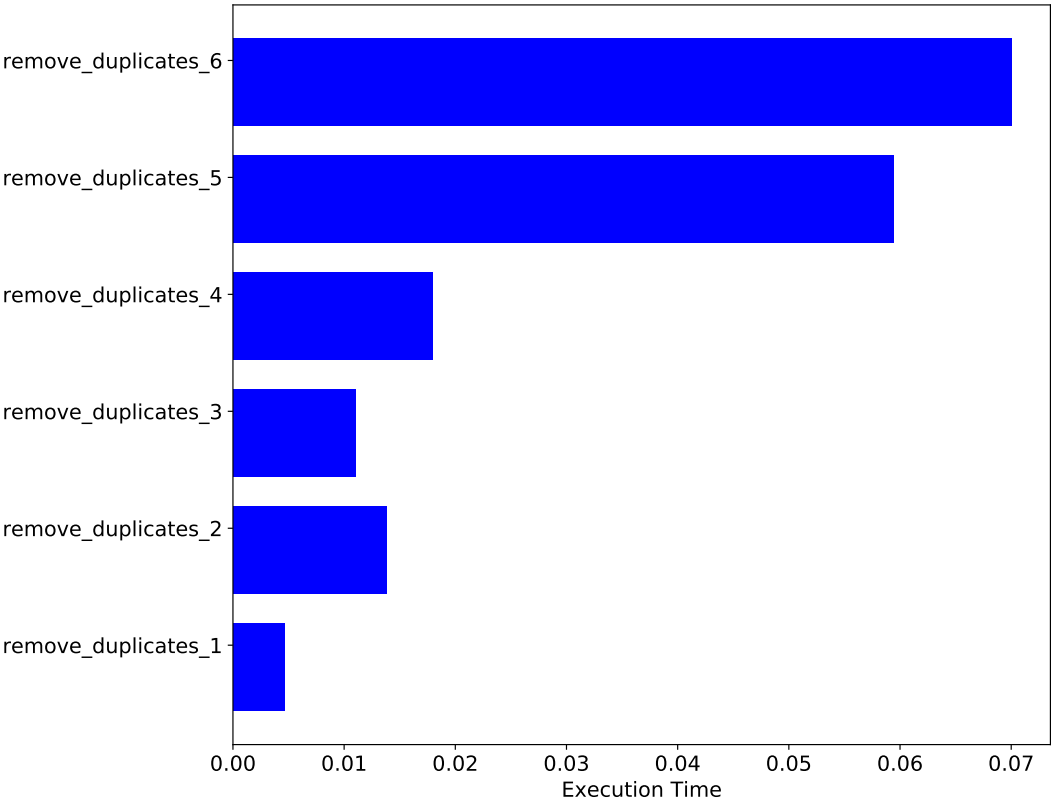
TIMES FOR THE FIRST 120000 NUMBERS



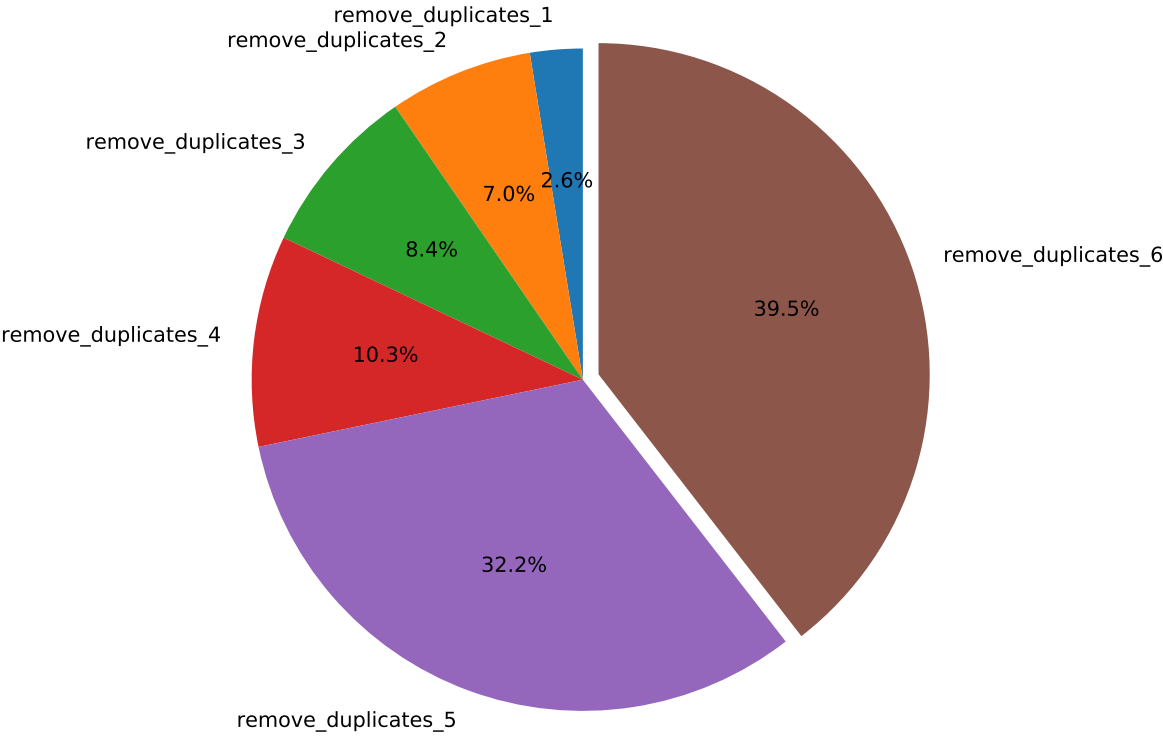
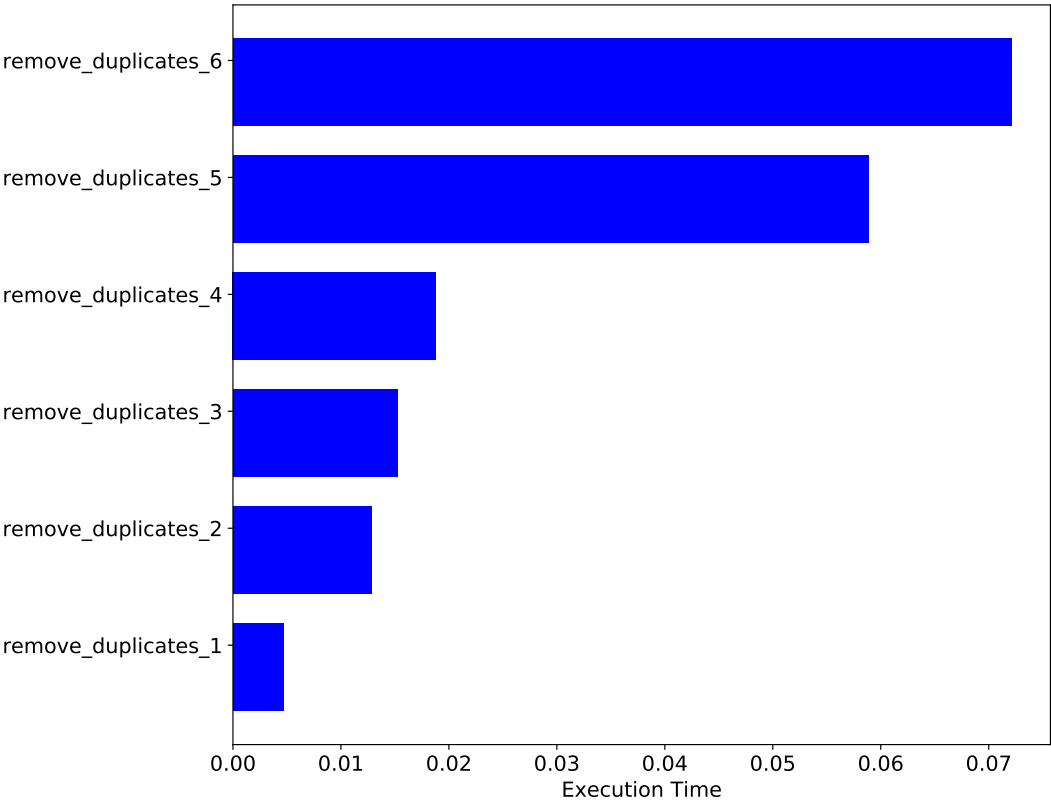
TIMES FOR THE FIRST 122000 NUMBERS



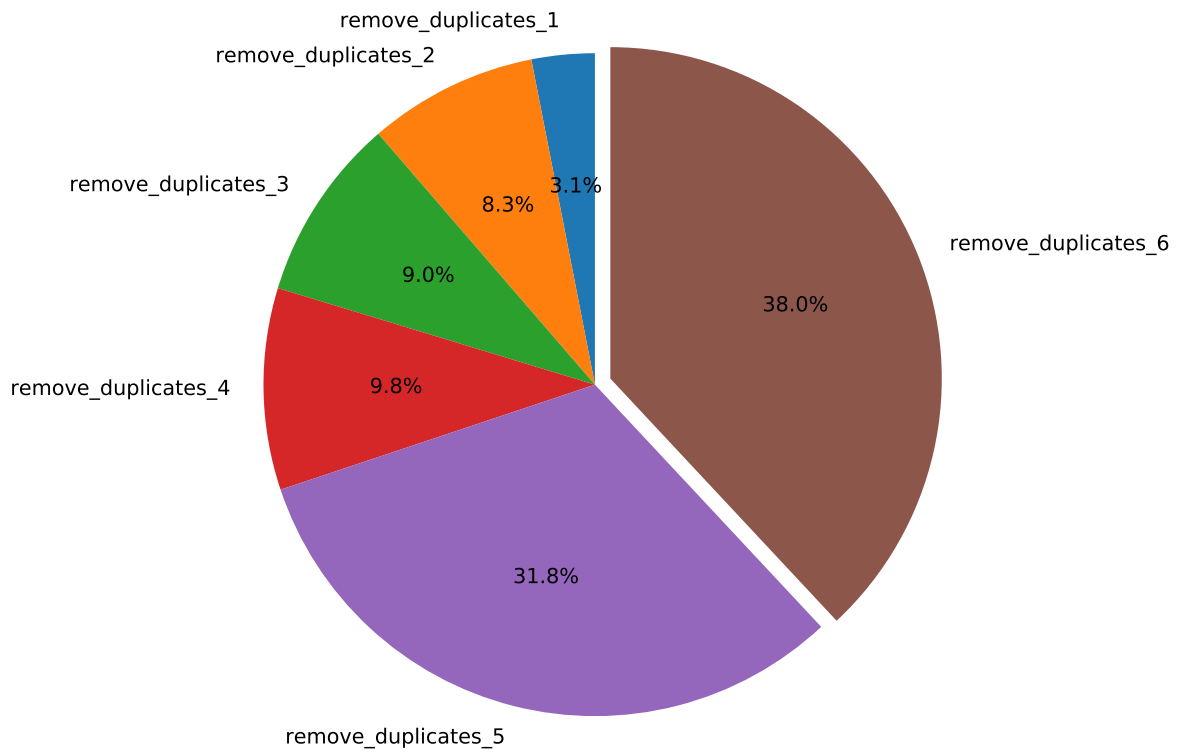
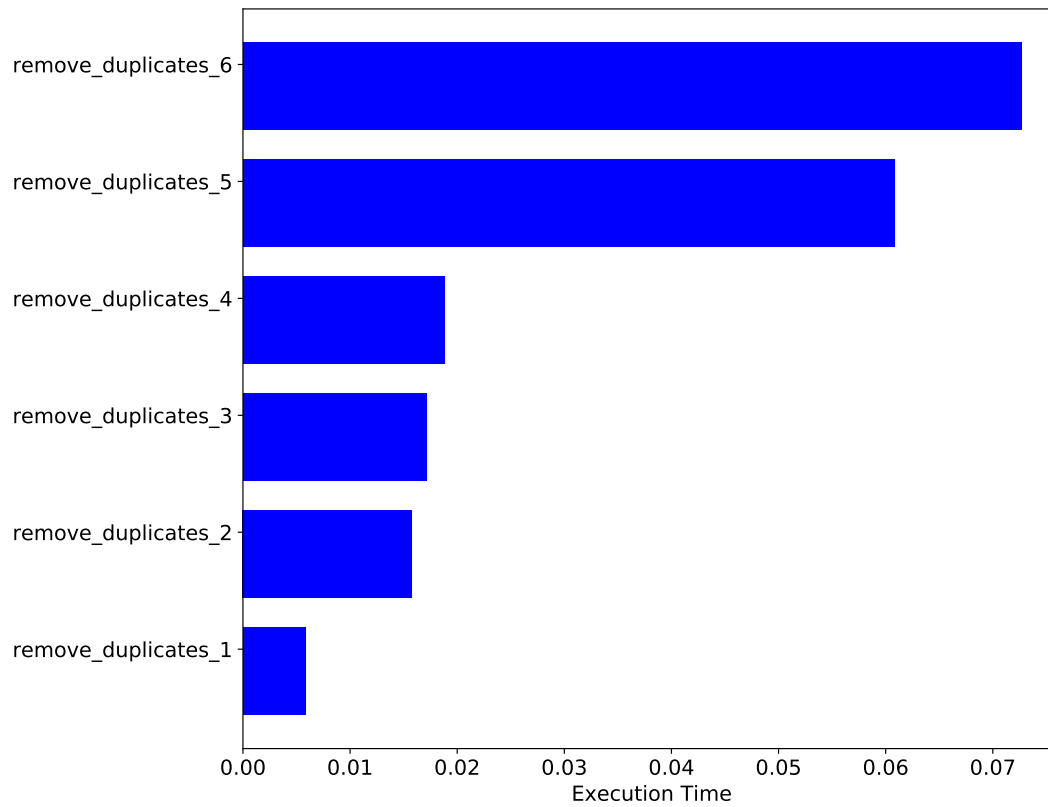
TIMES FOR THE FIRST 124000 NUMBERS



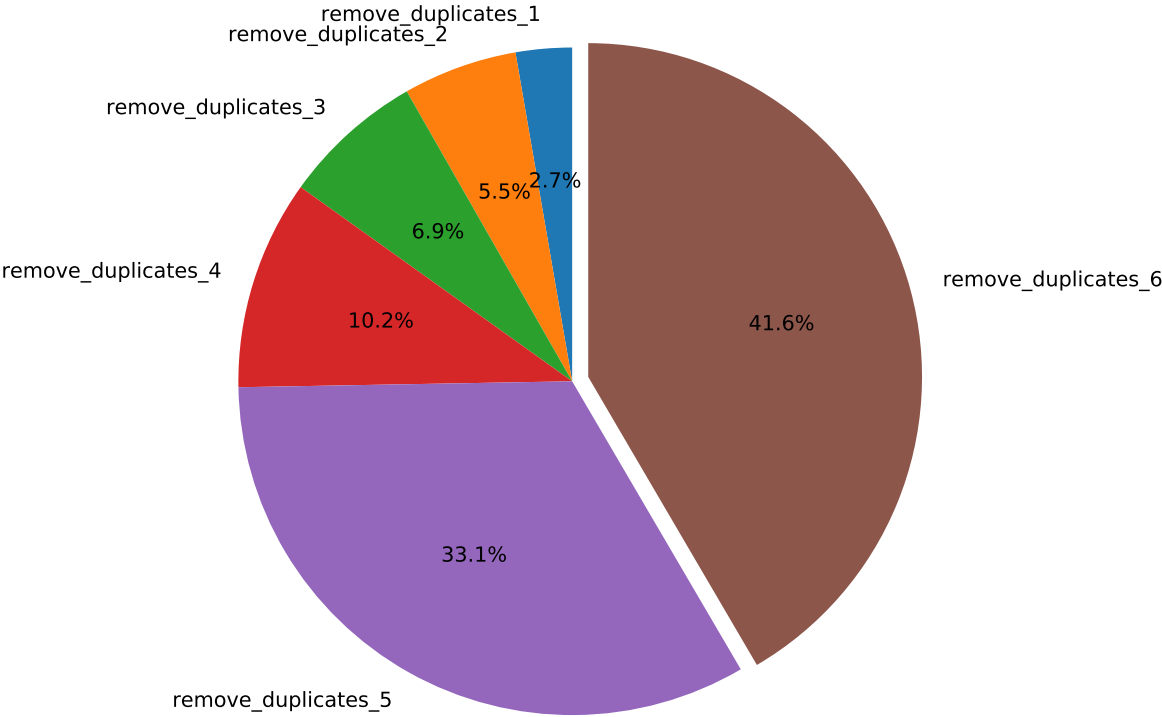
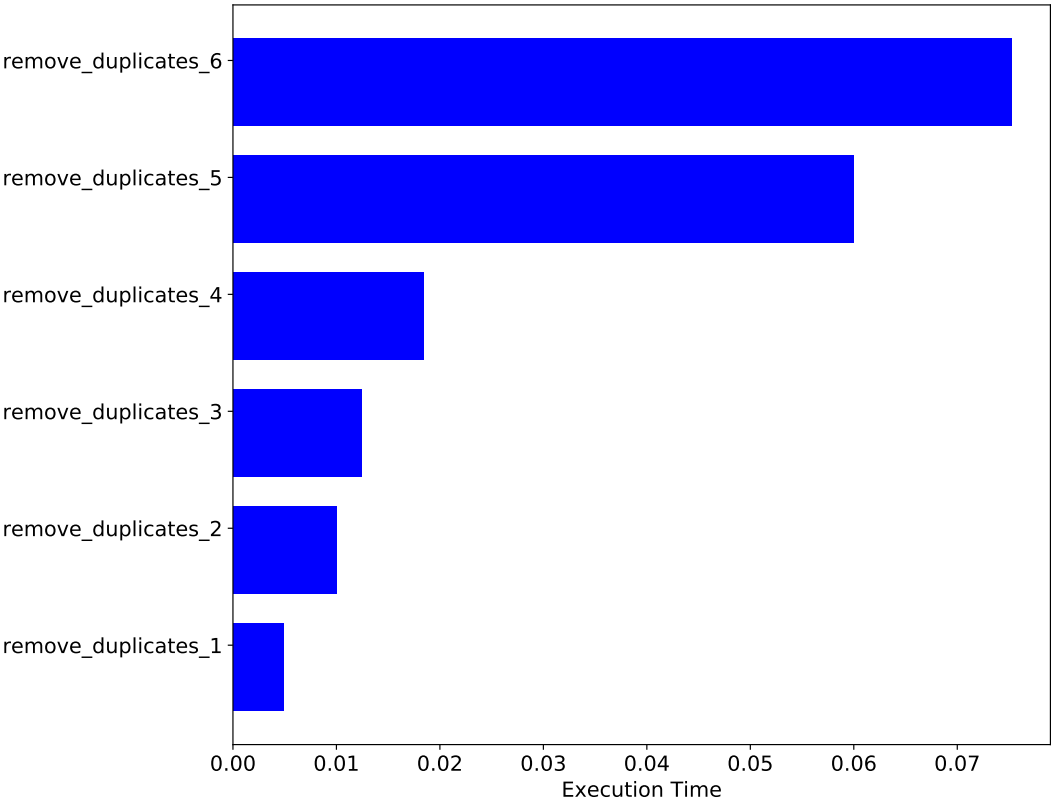
TIMES FOR THE FIRST 126000 NUMBERS



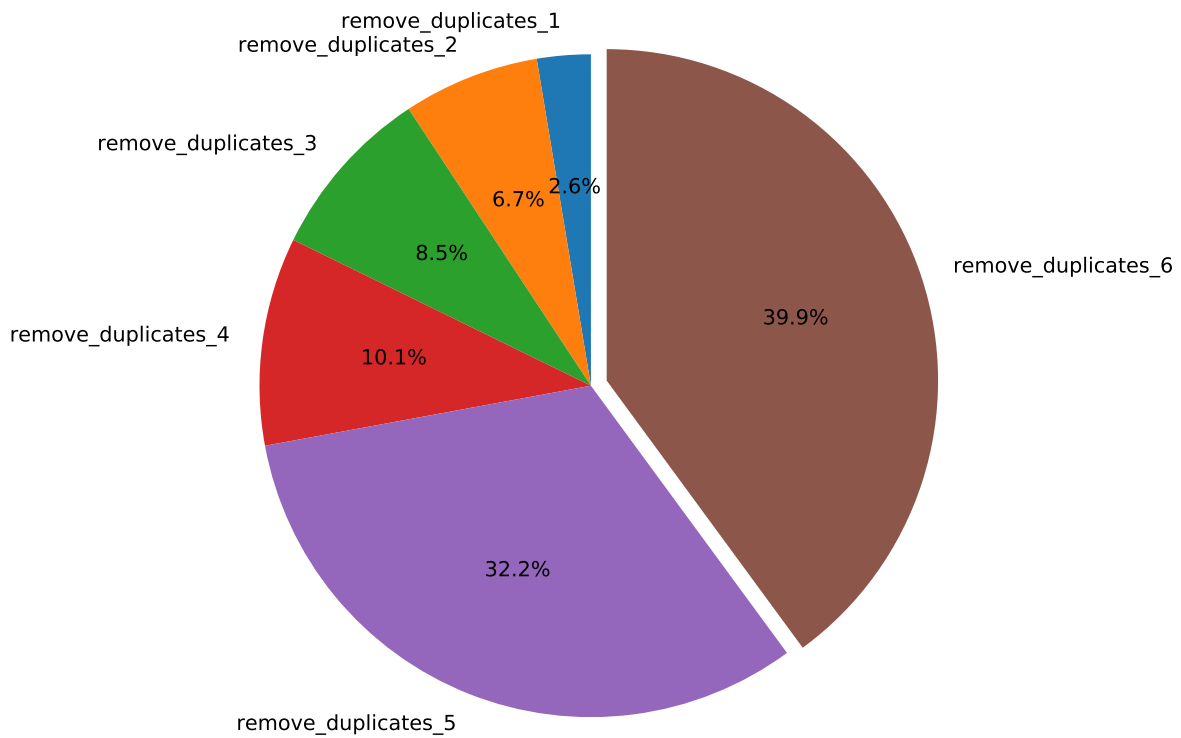
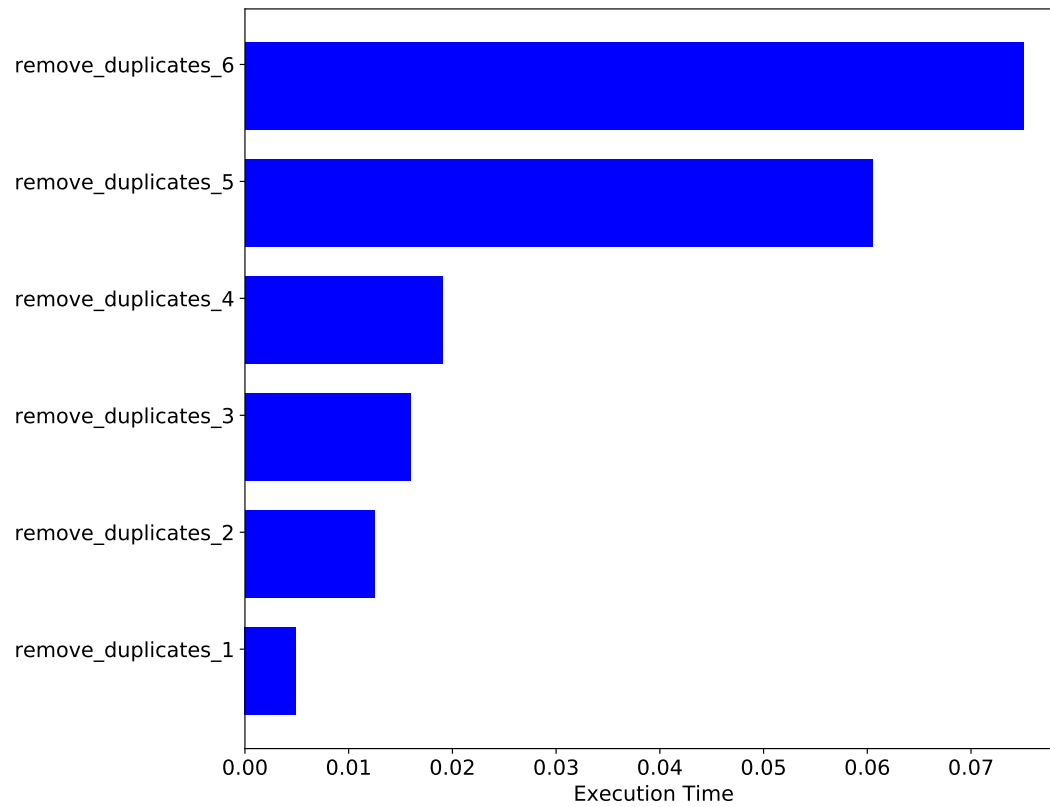
TIMES FOR THE FIRST 128000 NUMBERS



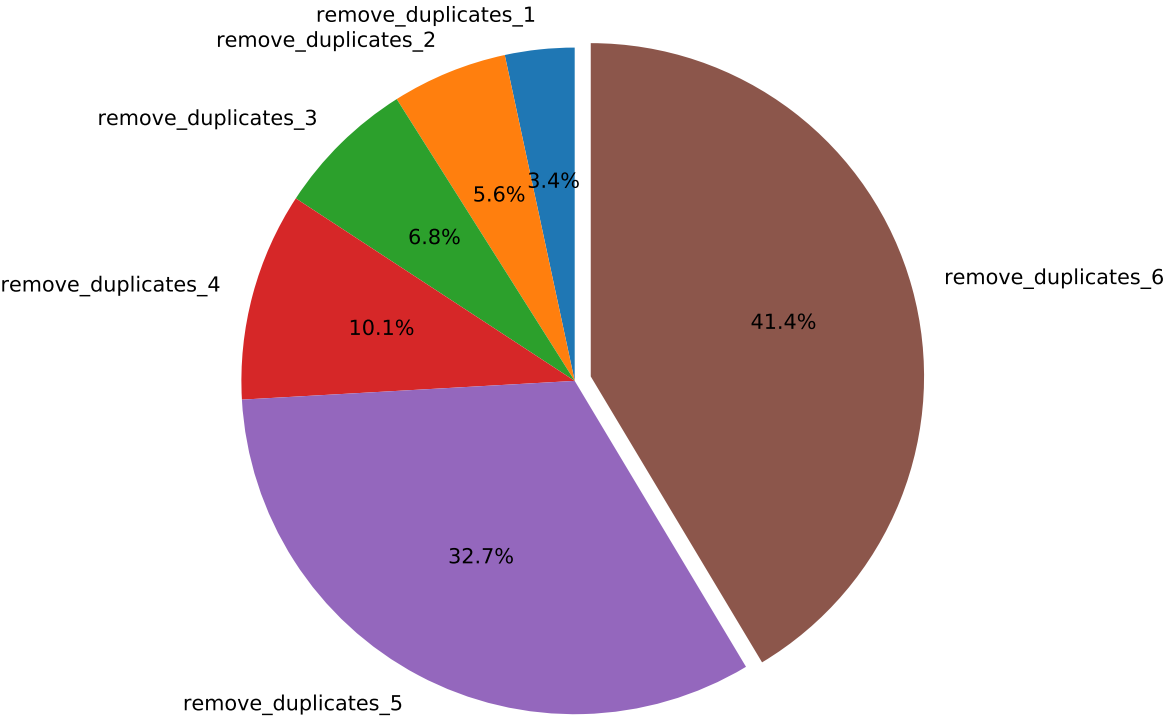
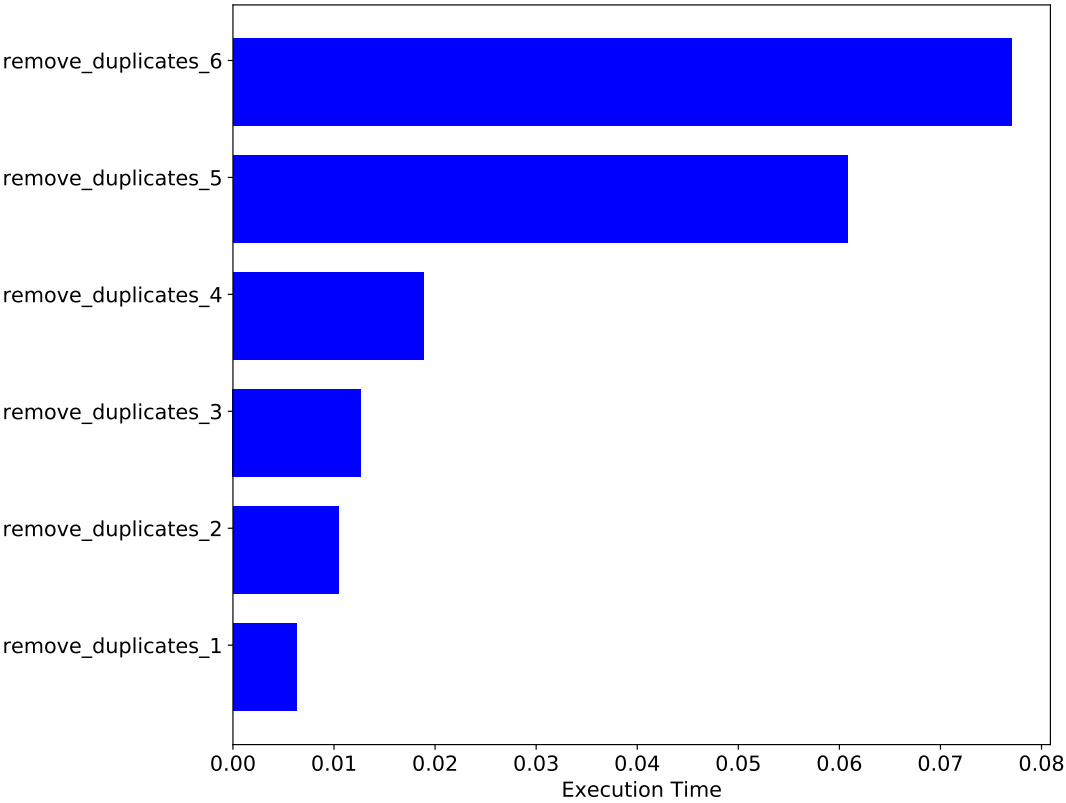
TIMES FOR THE FIRST 130000 NUMBERS



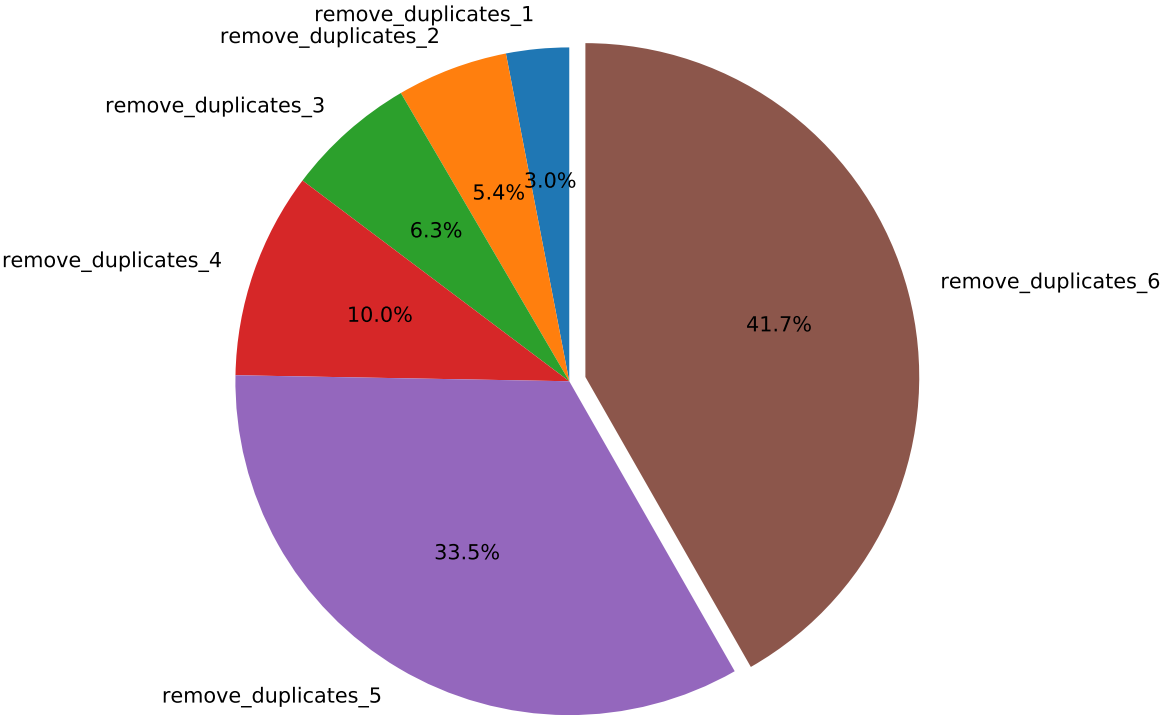
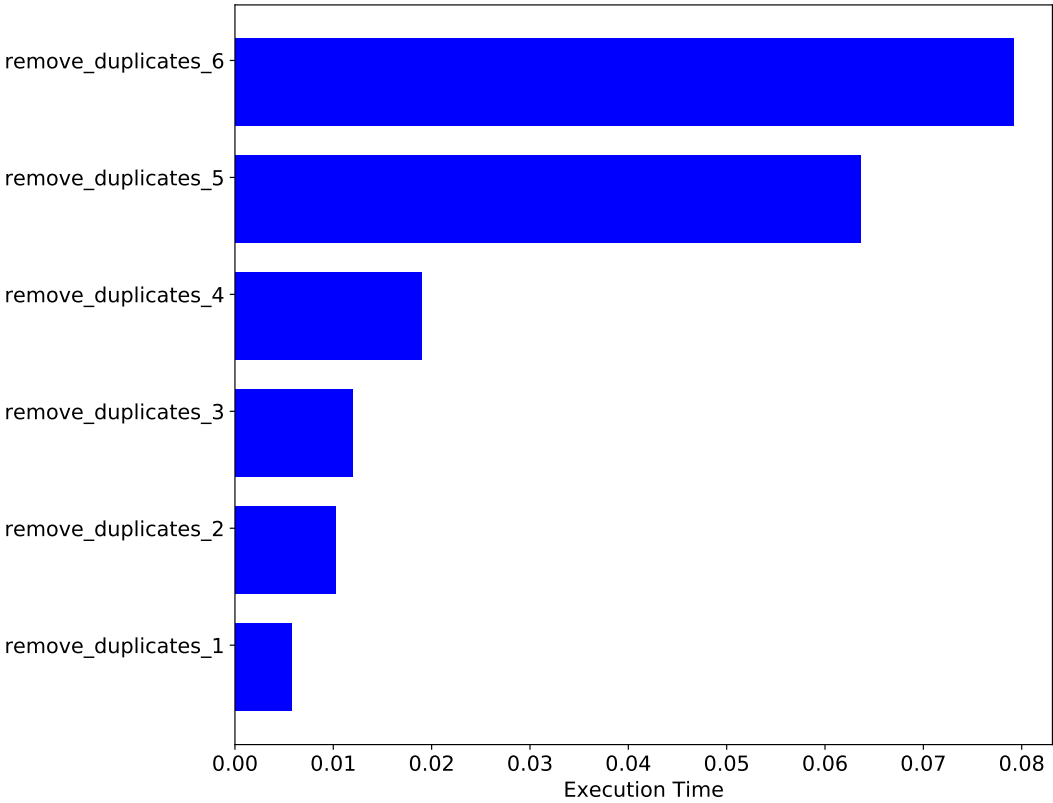
TIMES FOR THE FIRST 132000 NUMBERS



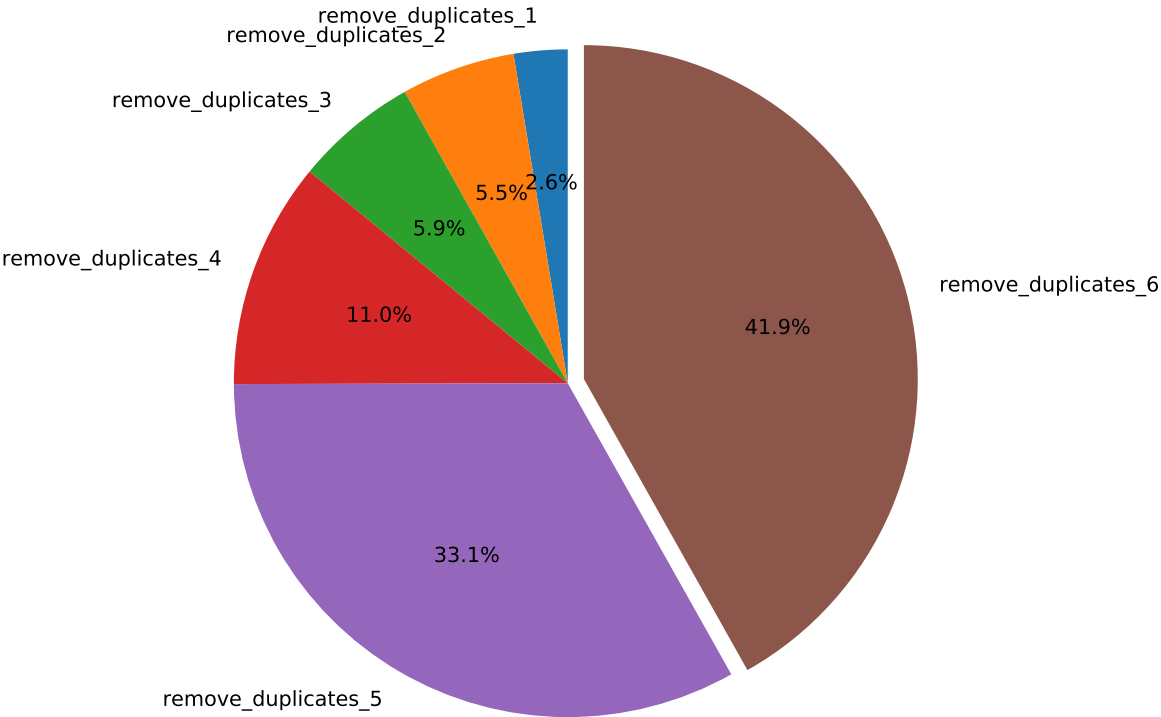
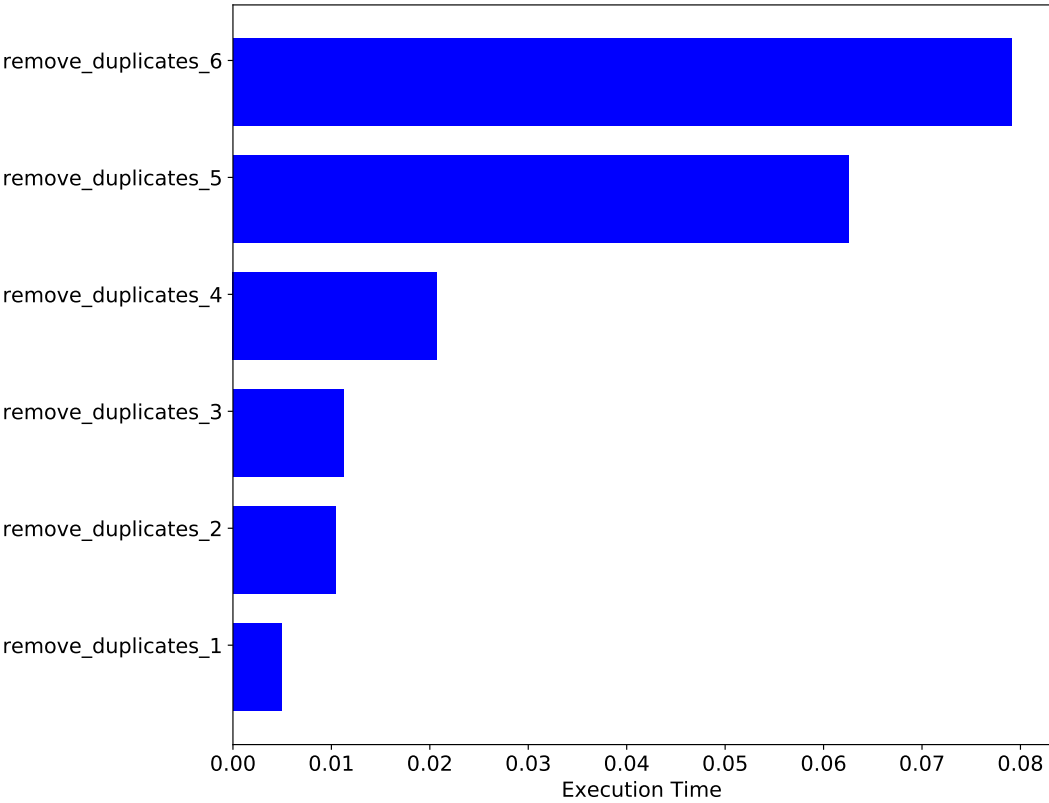
TIMES FOR THE FIRST 134000 NUMBERS



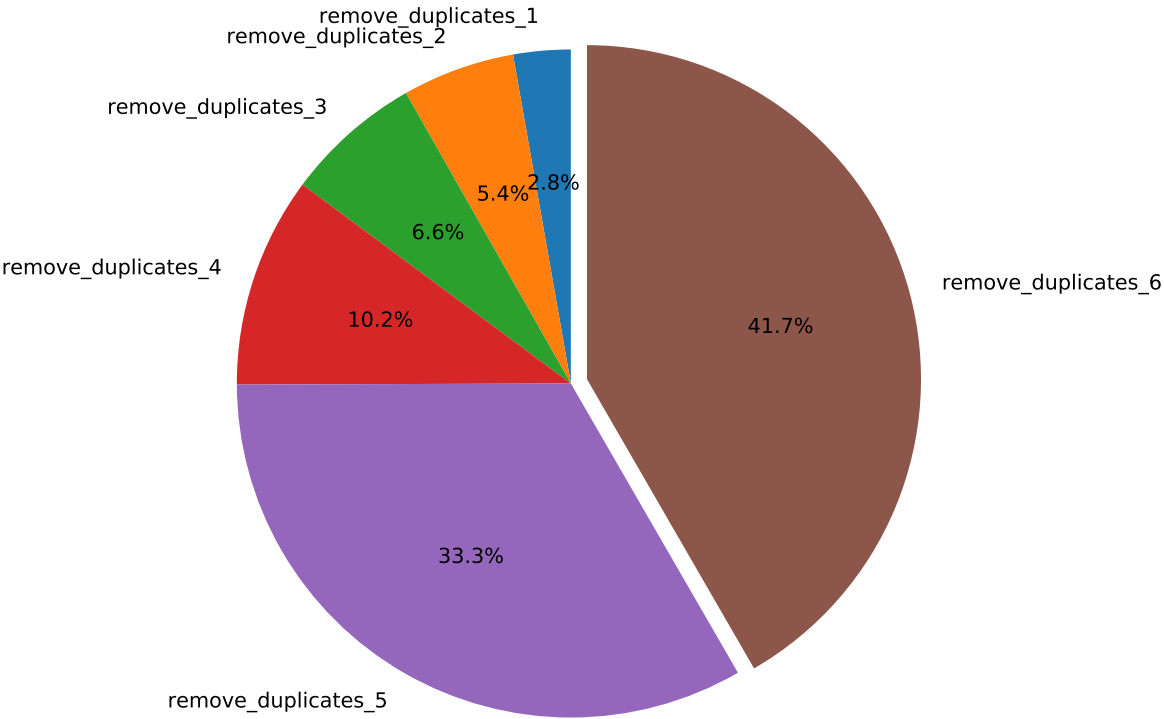
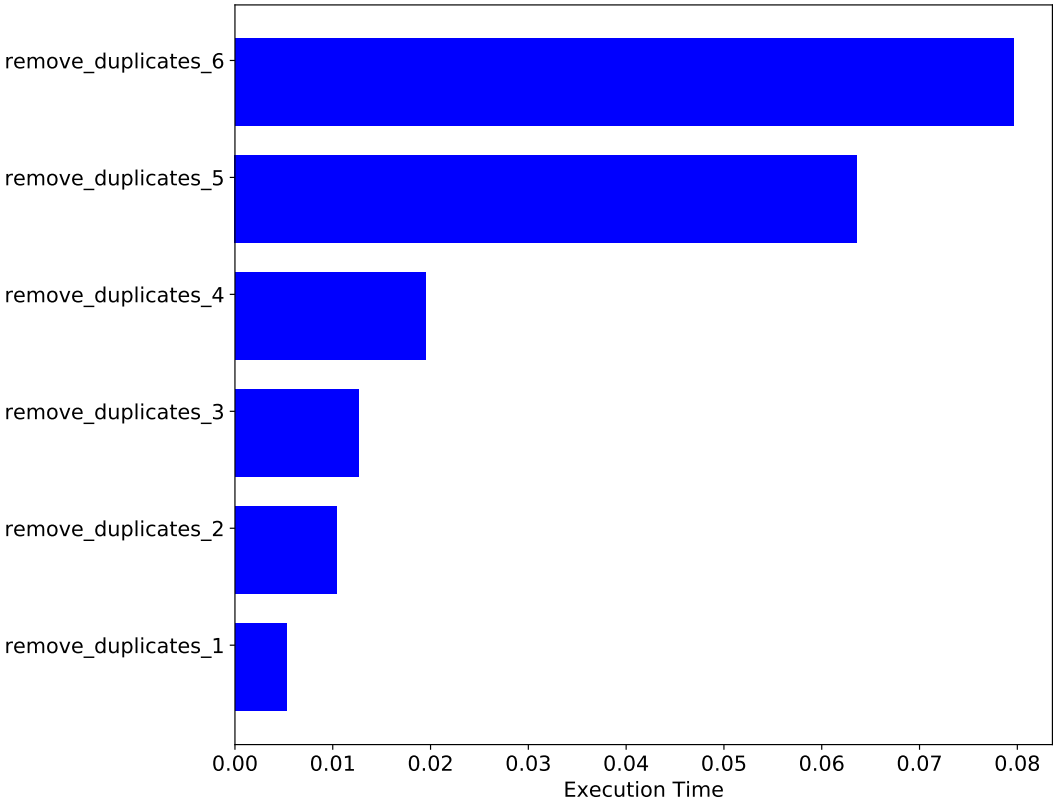
TIMES FOR THE FIRST 136000 NUMBERS



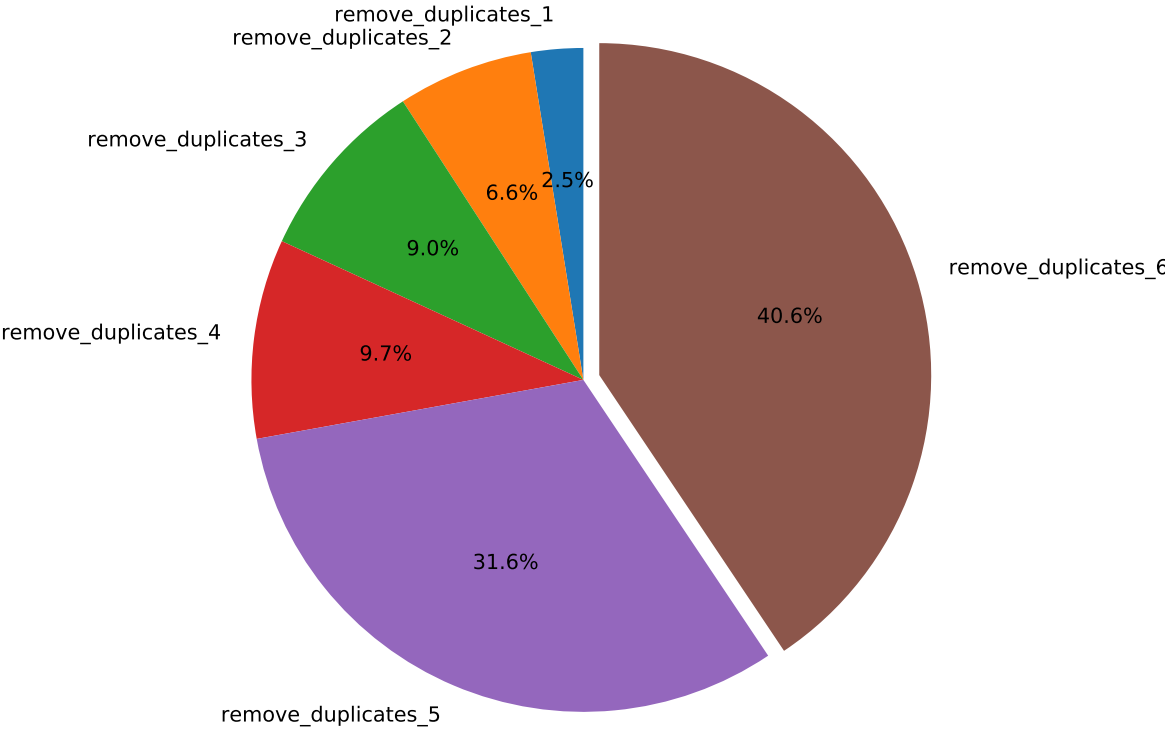
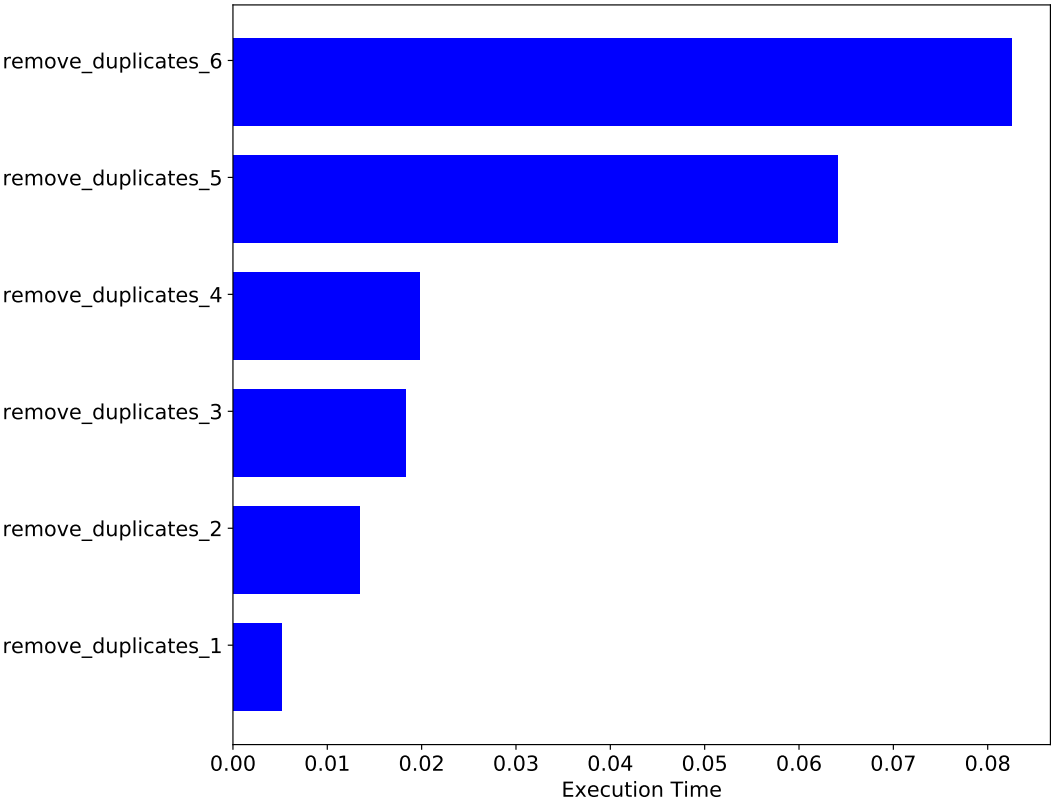
TIMES FOR THE FIRST 138000 NUMBERS



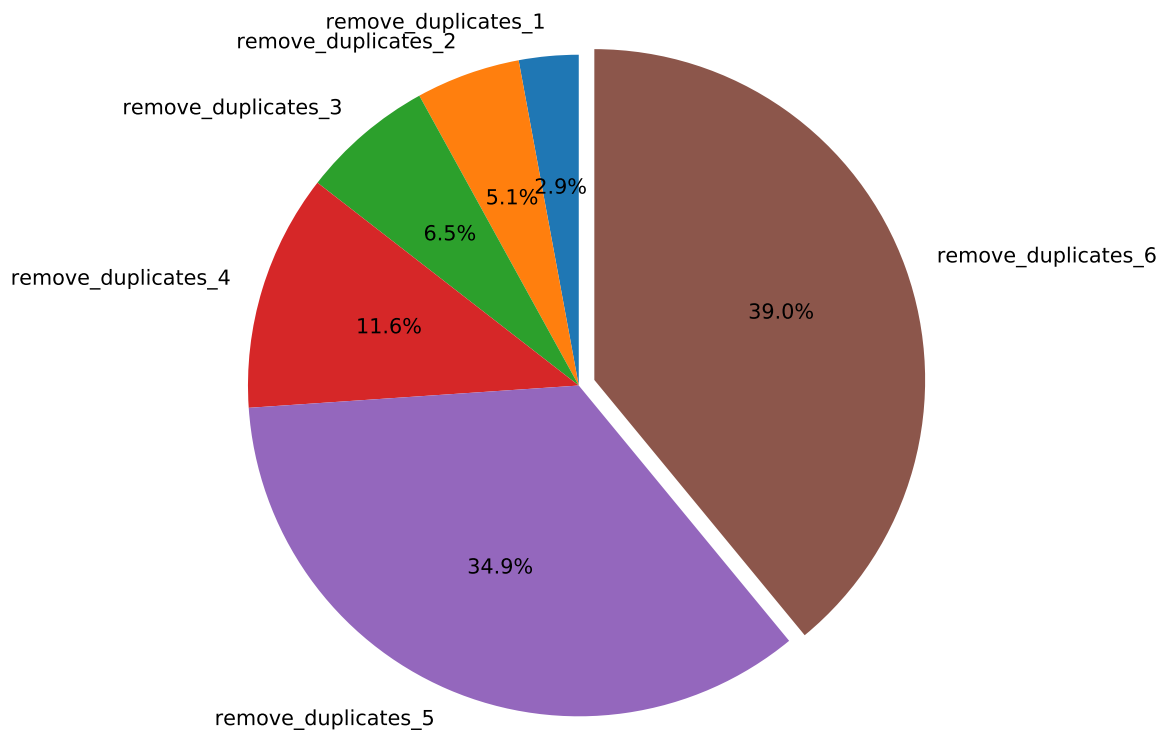
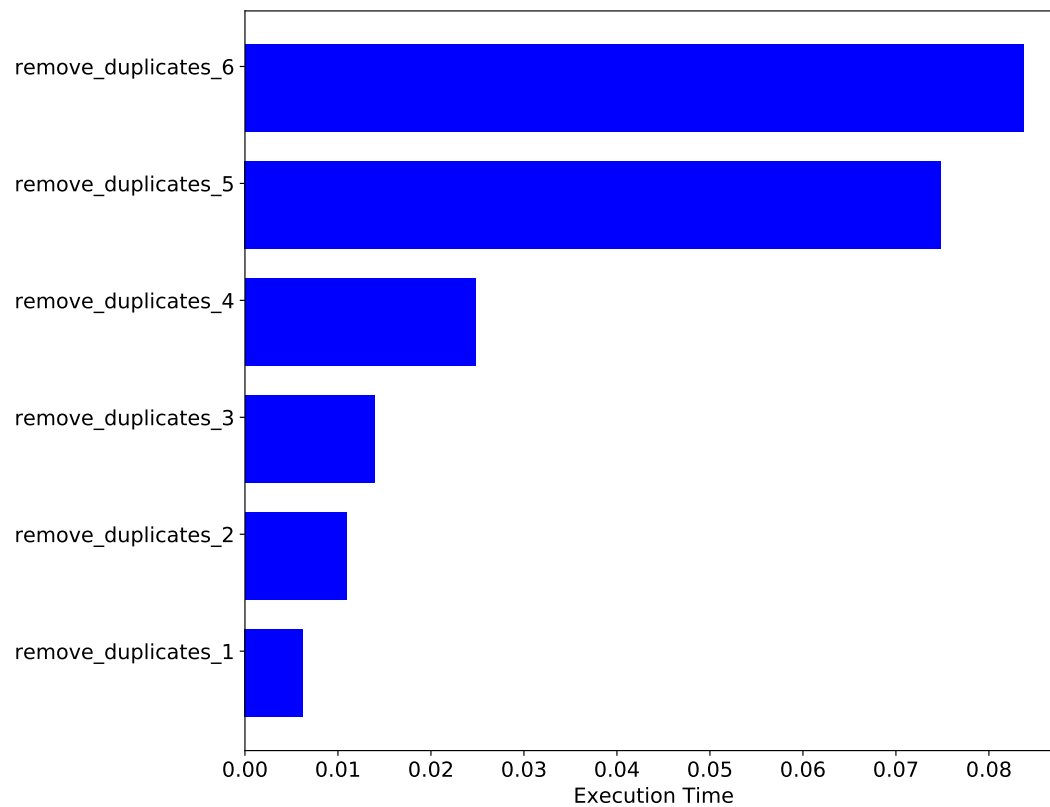
TIMES FOR THE FIRST 140000 NUMBERS



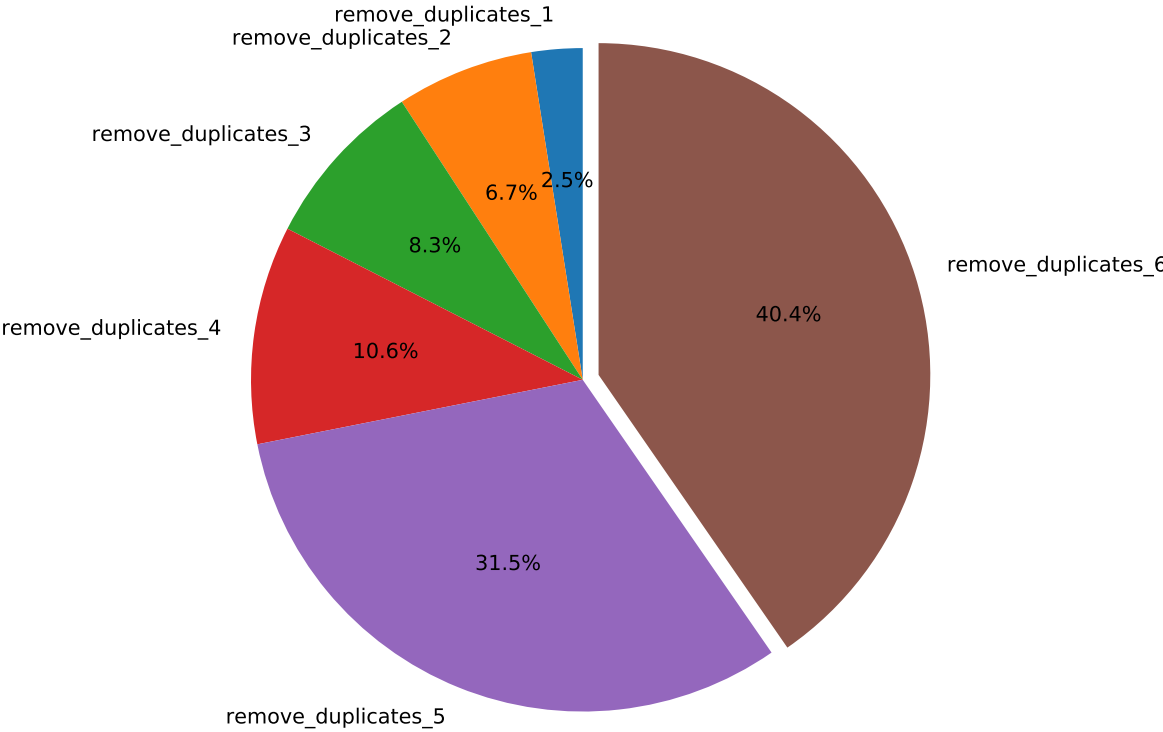
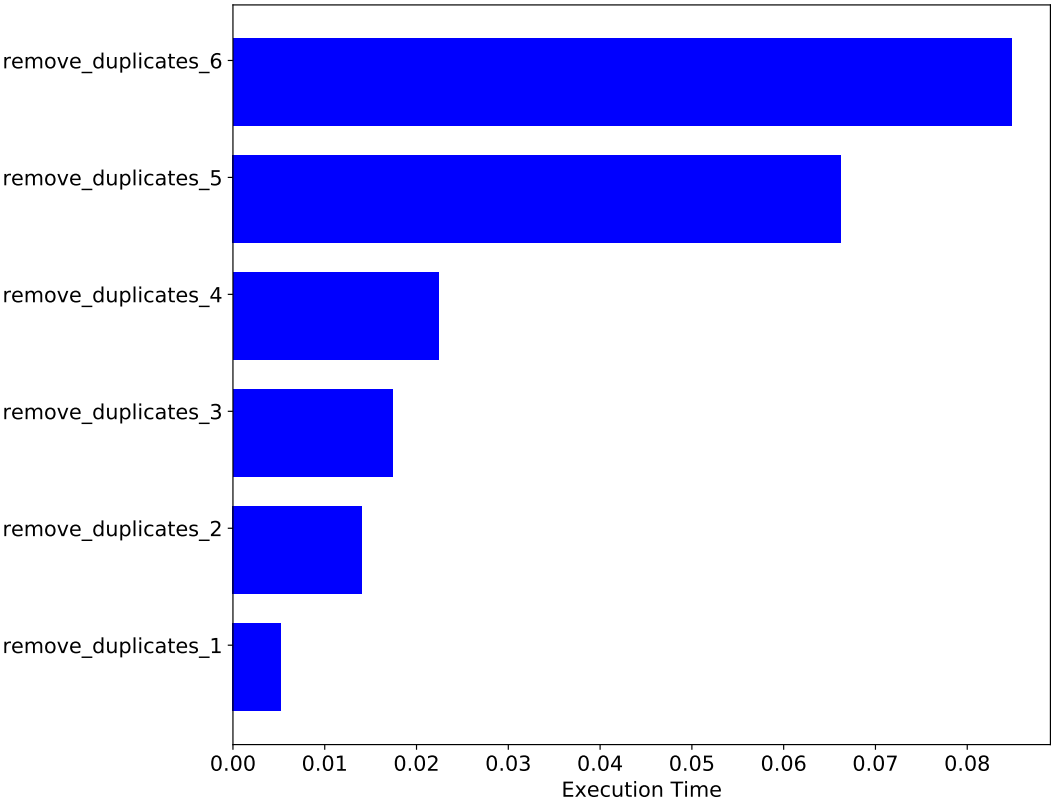
TIMES FOR THE FIRST 142000 NUMBERS



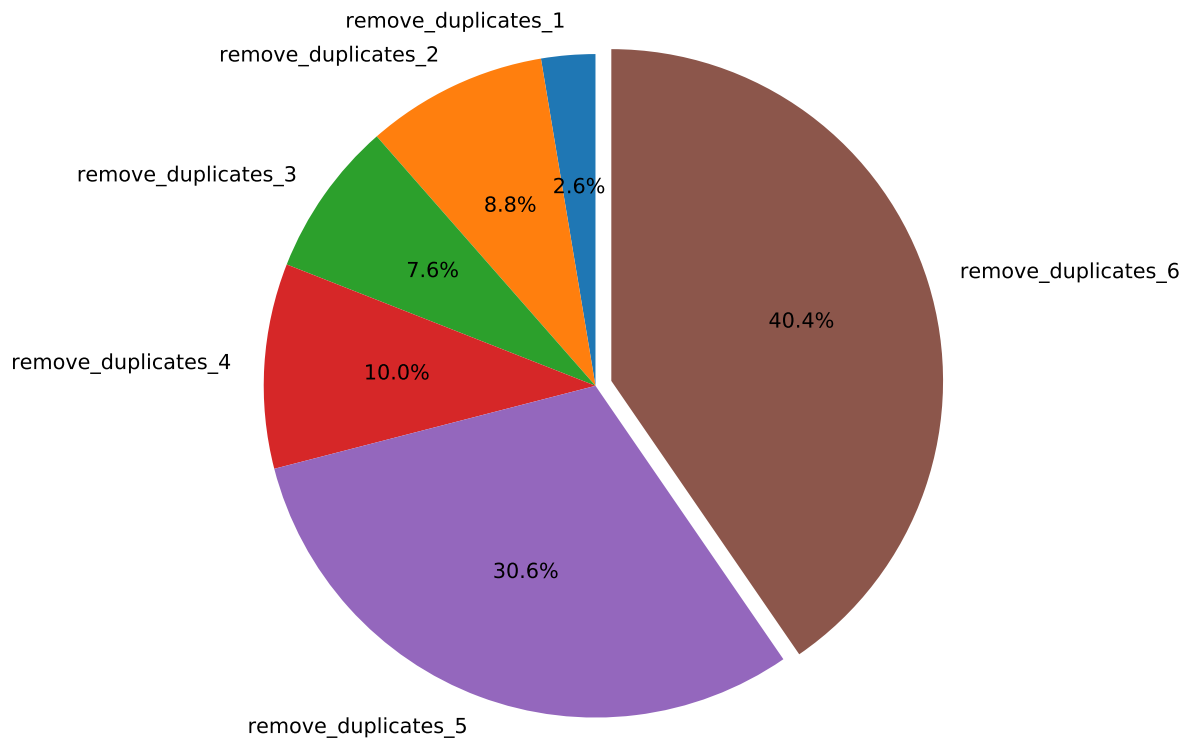
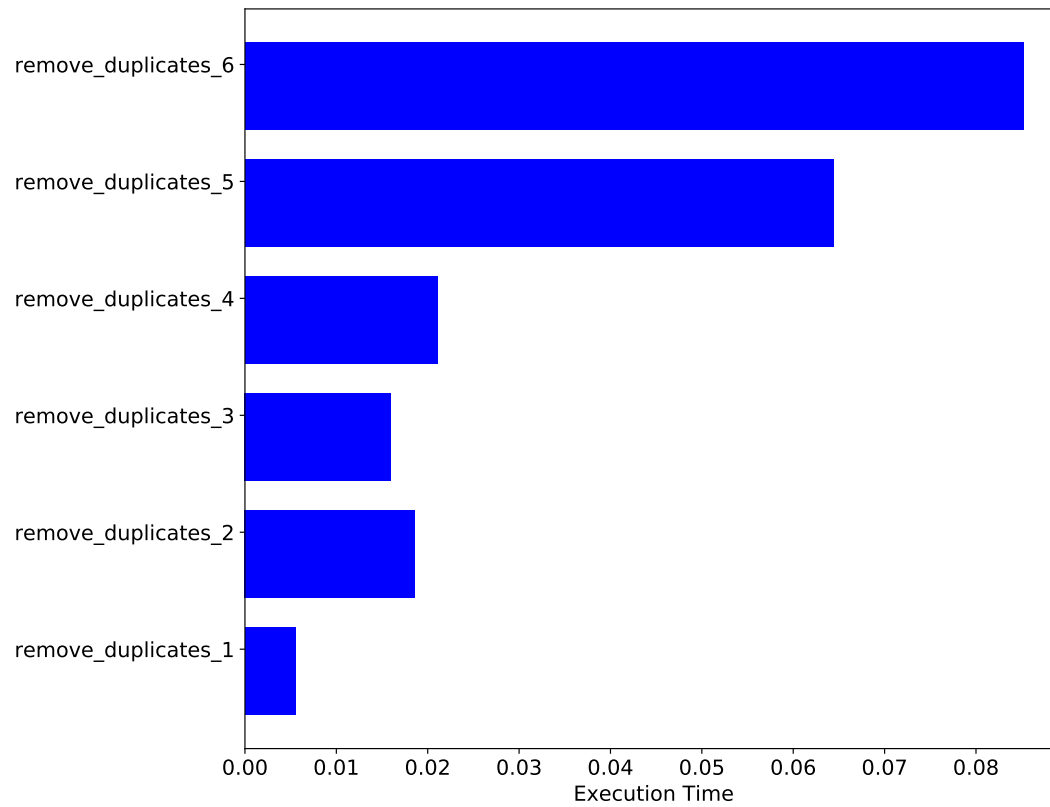
TIMES FOR THE FIRST 144000 NUMBERS



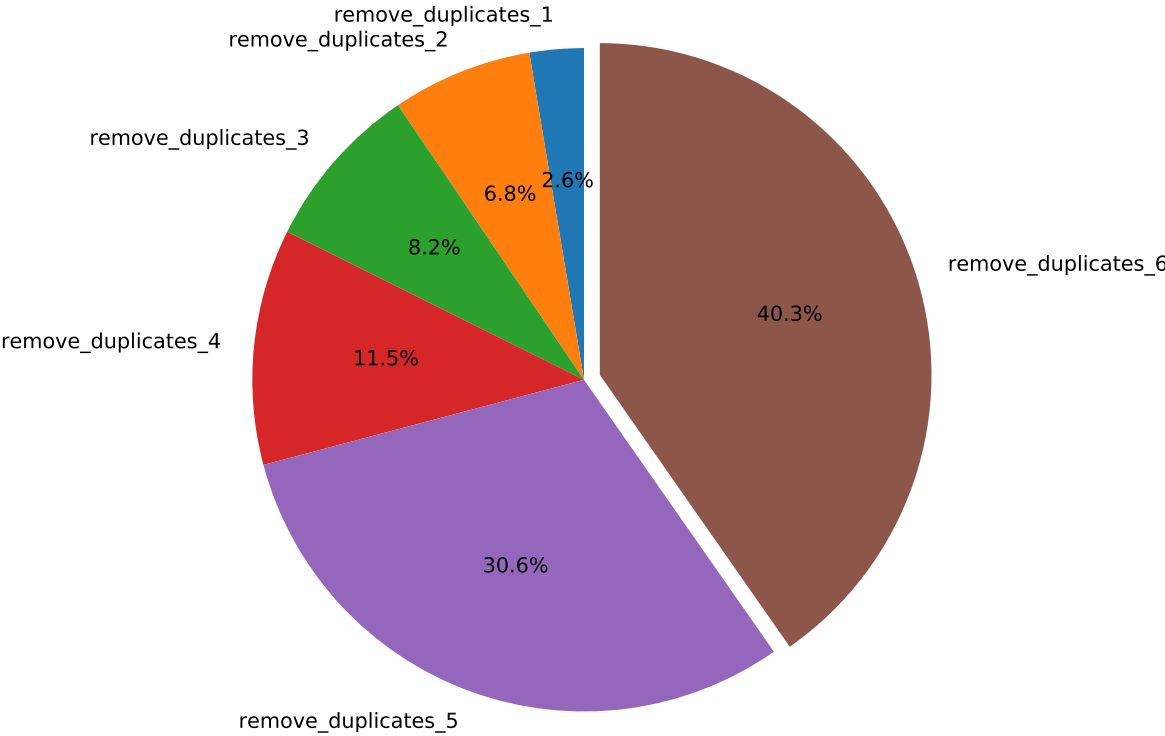
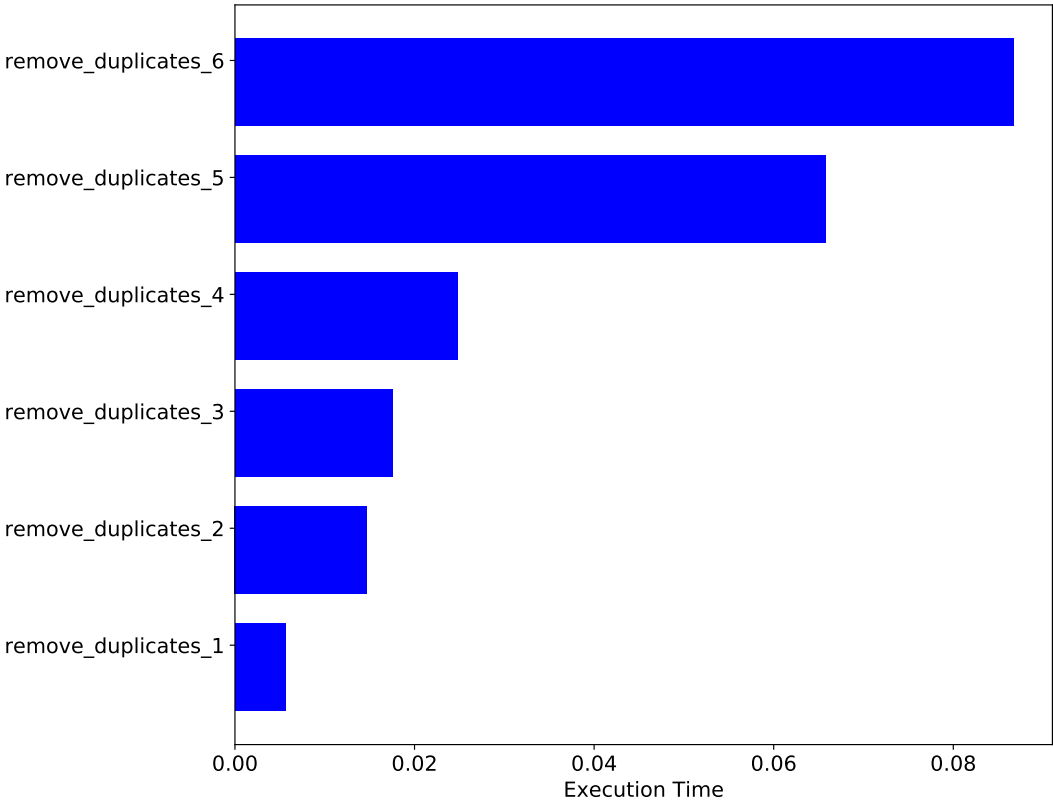
TIMES FOR THE FIRST 146000 NUMBERS



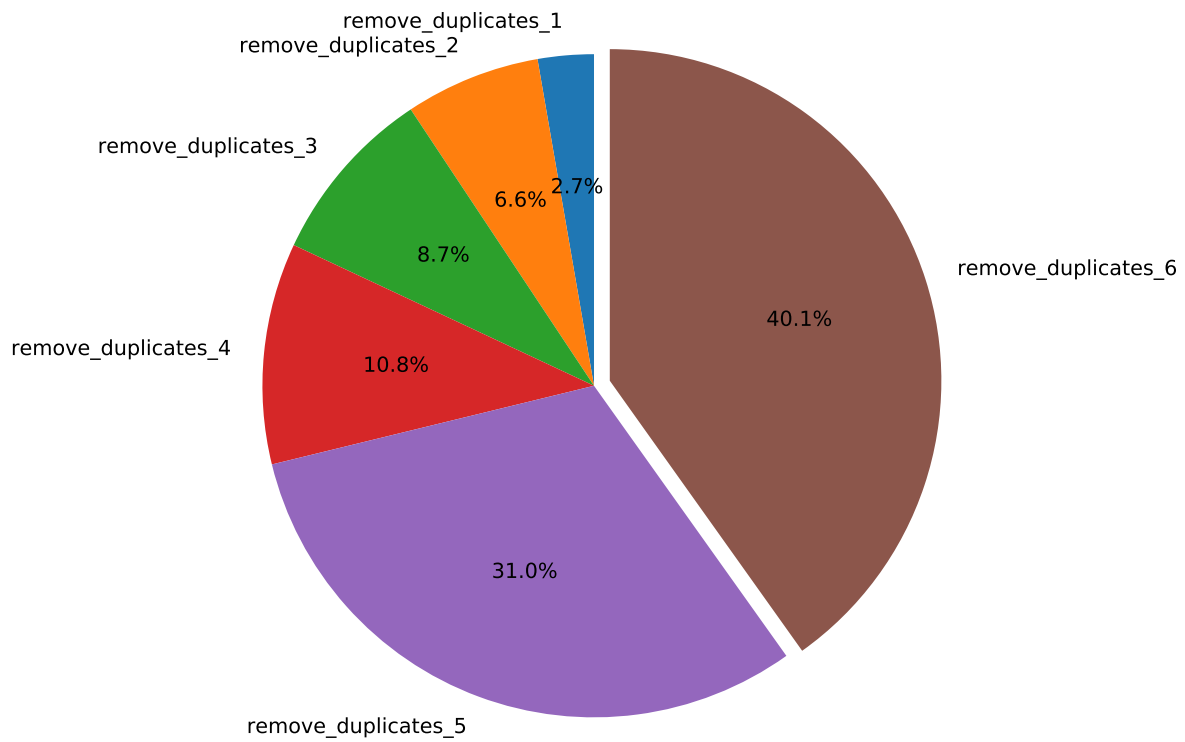
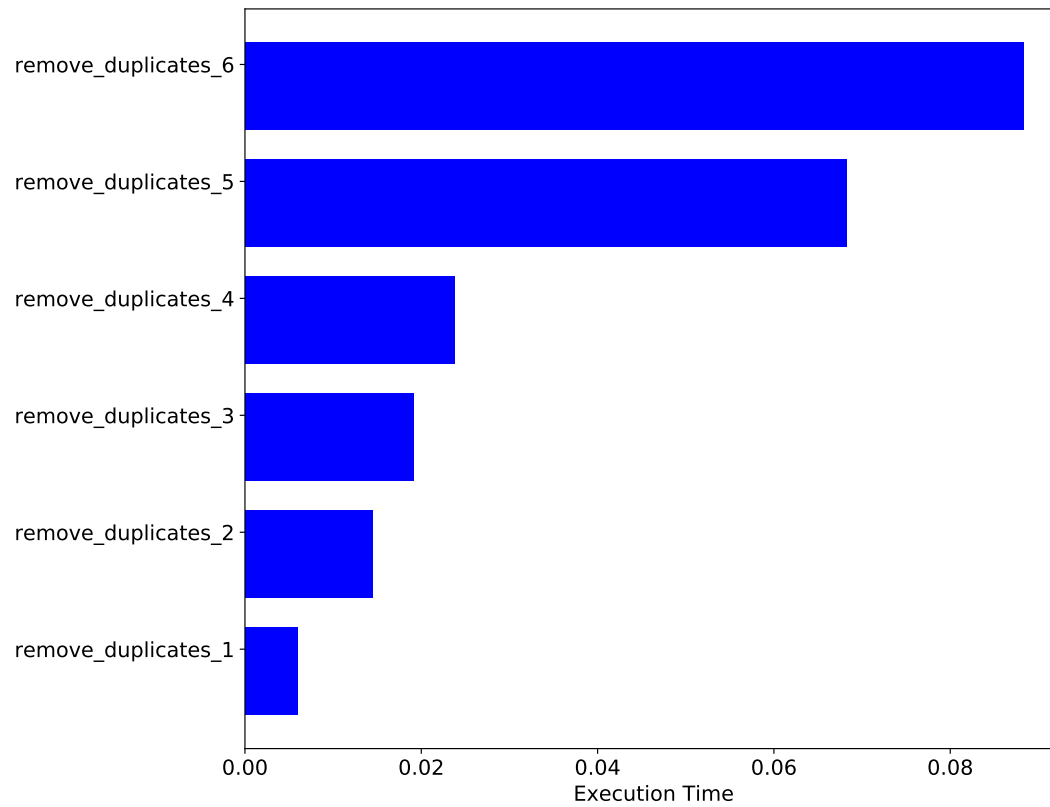
TIMES FOR THE FIRST 148000 NUMBERS



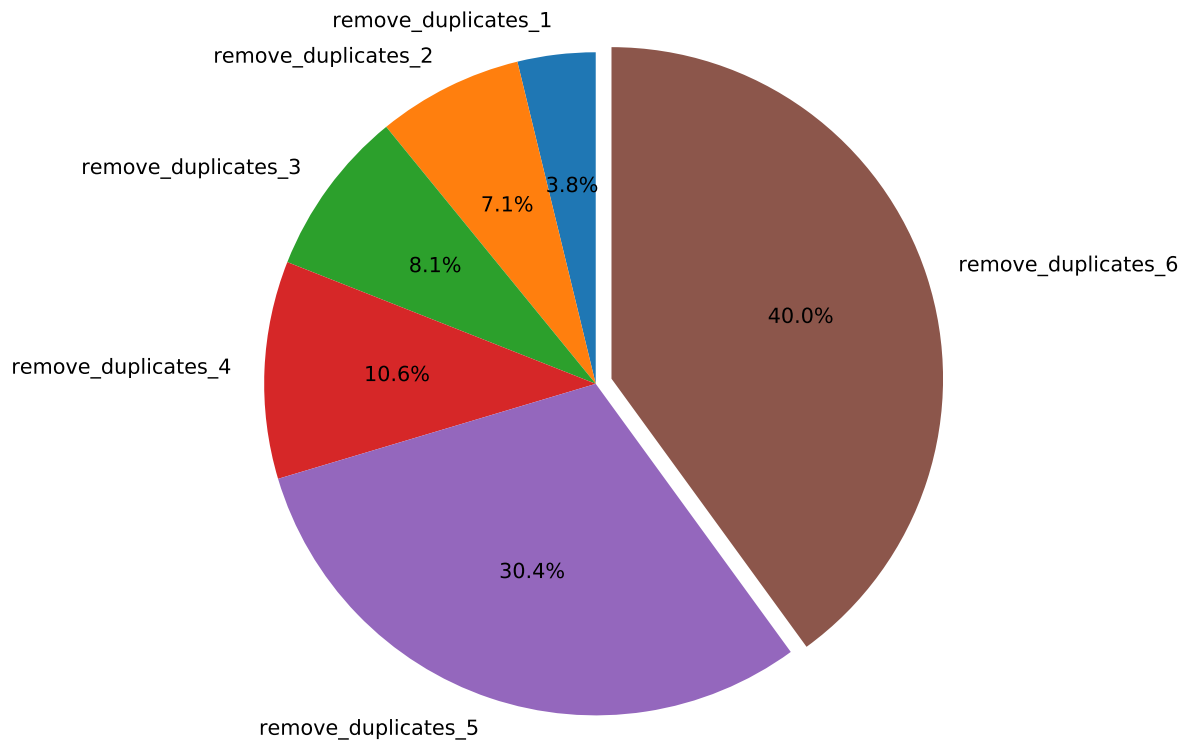
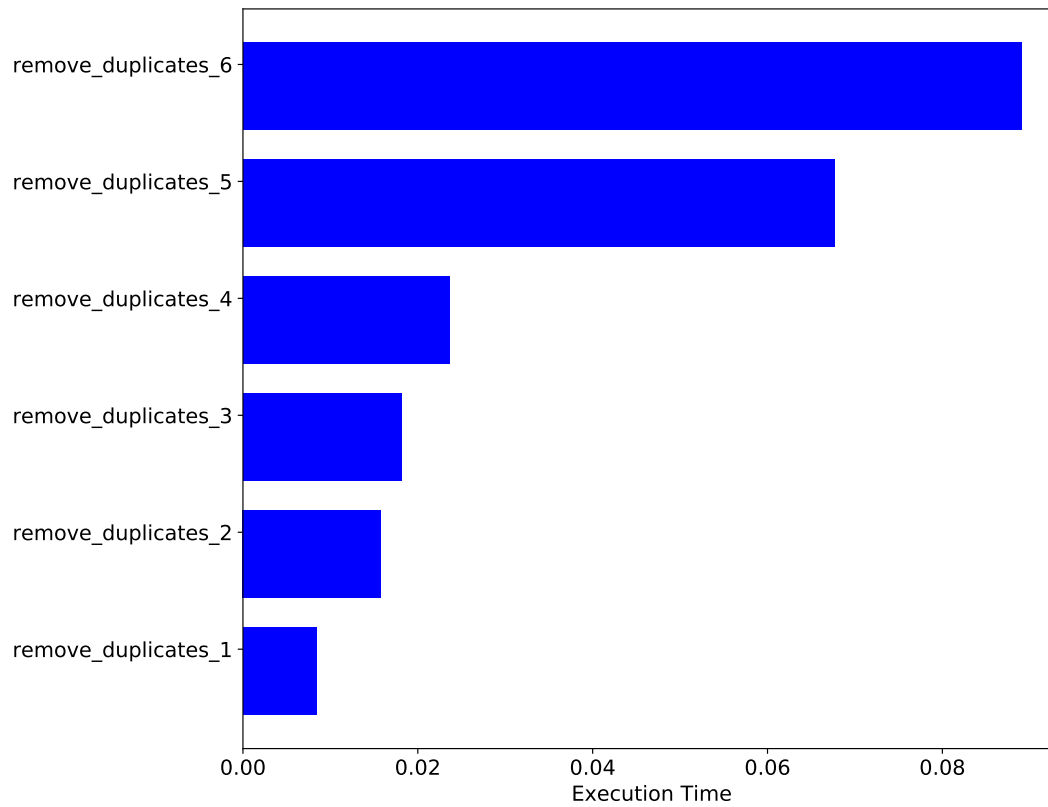
TIMES FOR THE FIRST 150000 NUMBERS



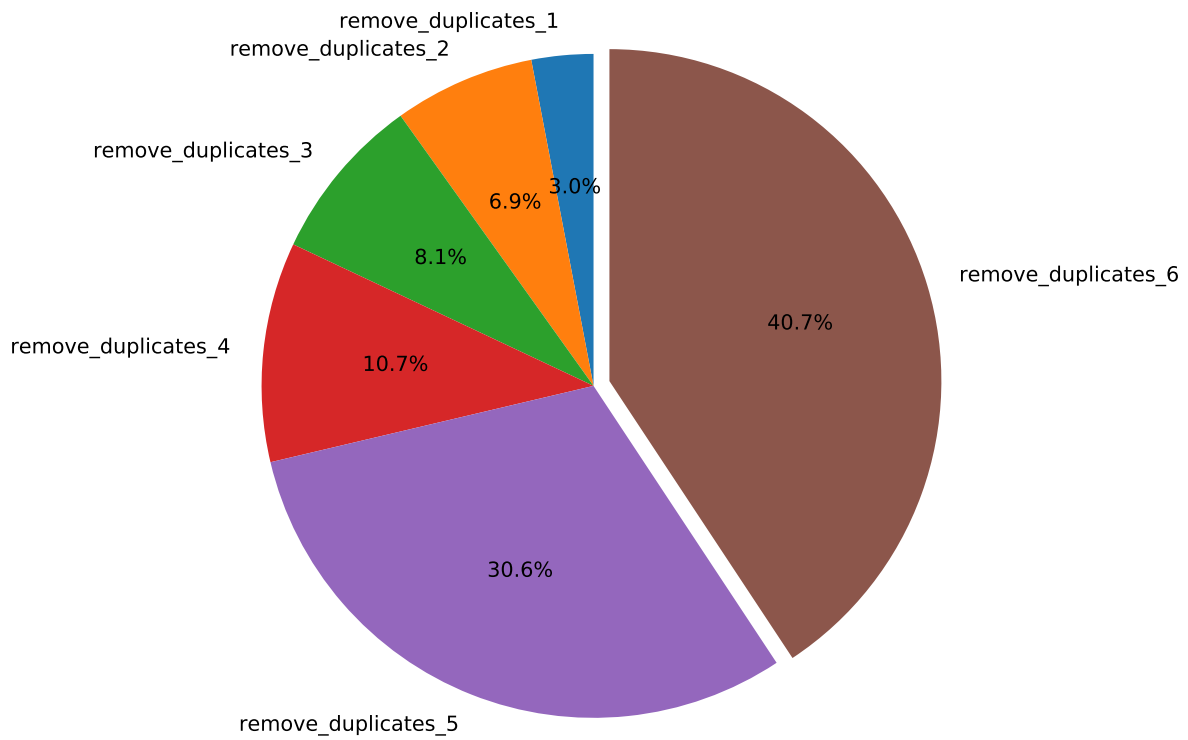
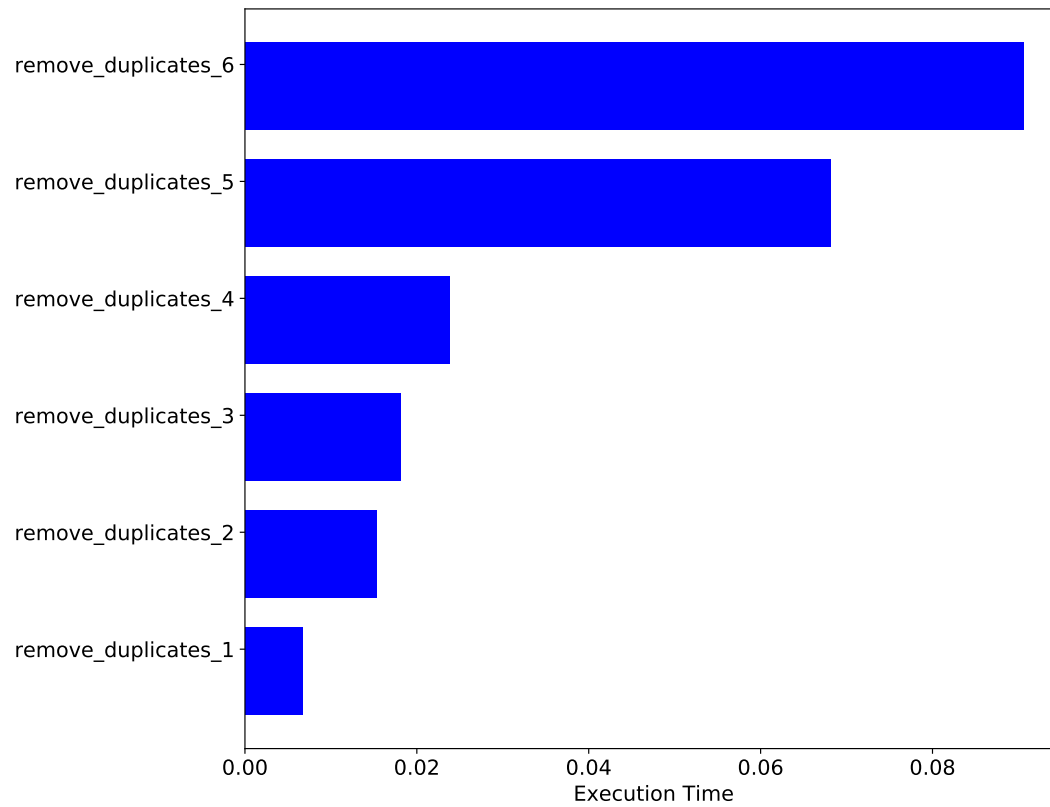
TIMES FOR THE FIRST 152000 NUMBERS



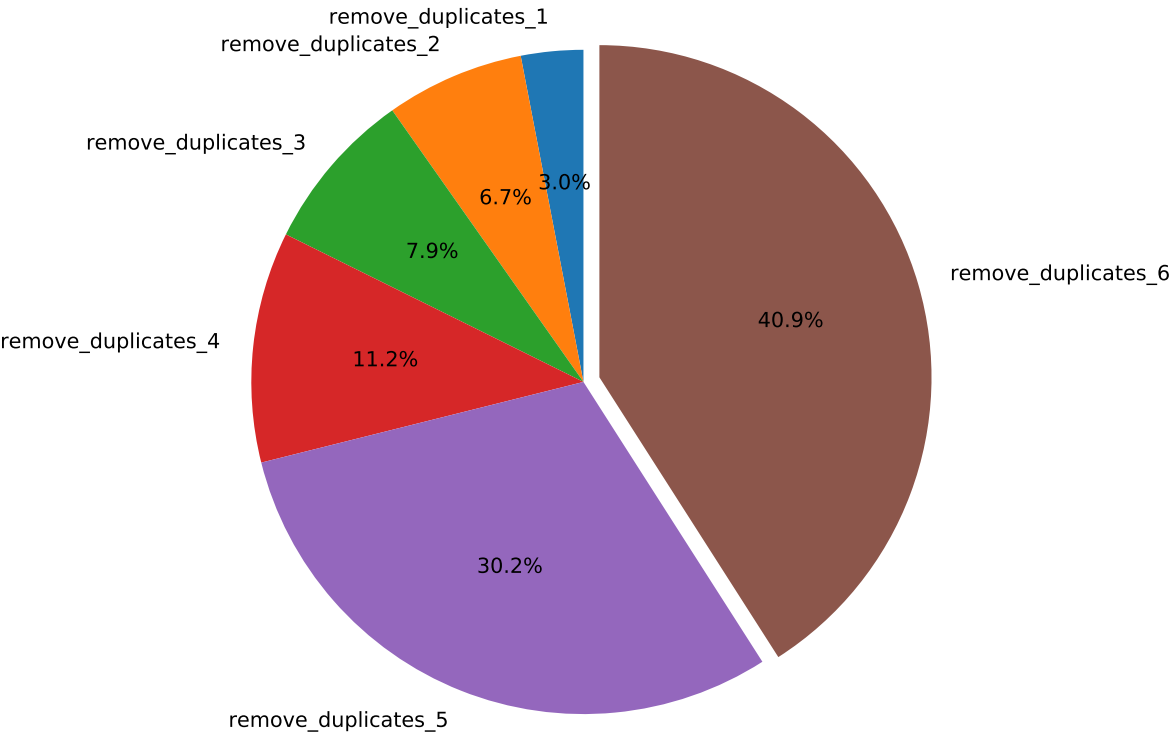
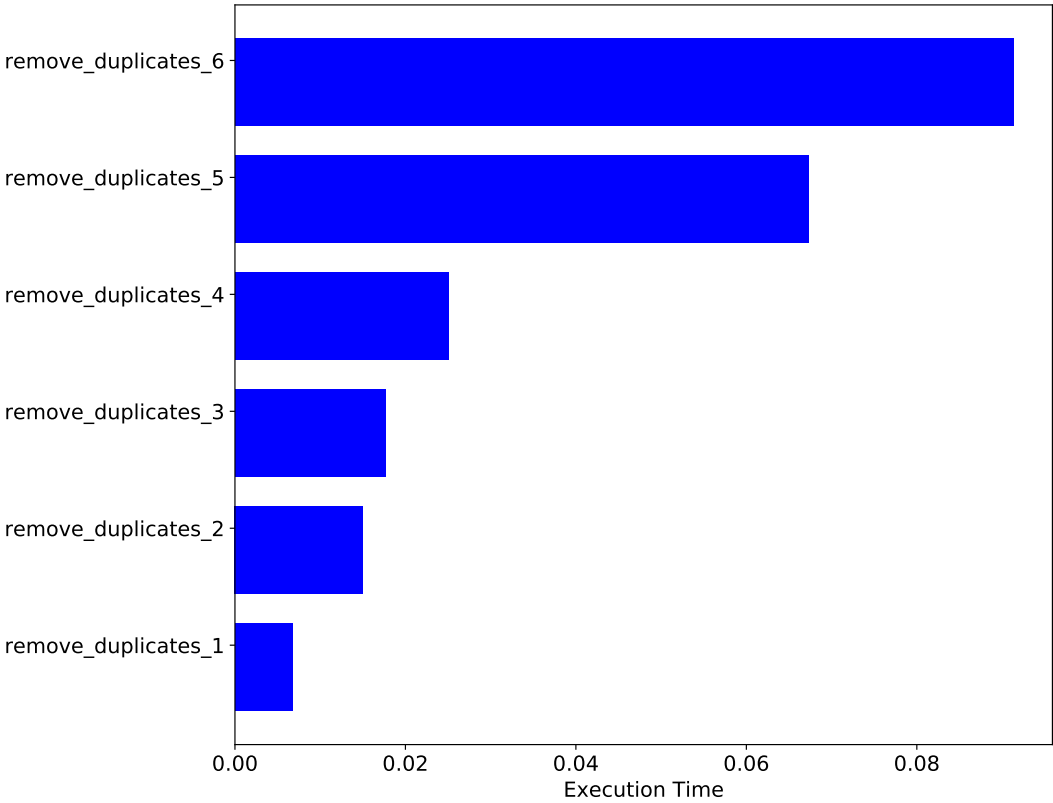
TIMES FOR THE FIRST 154000 NUMBERS



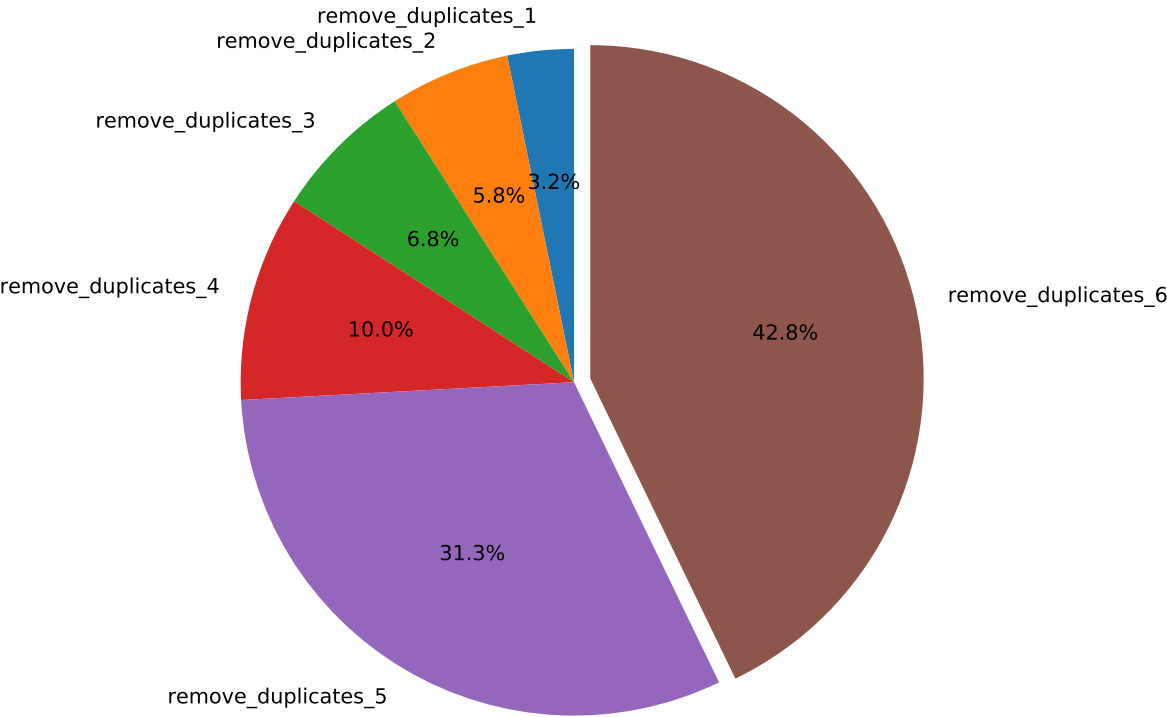
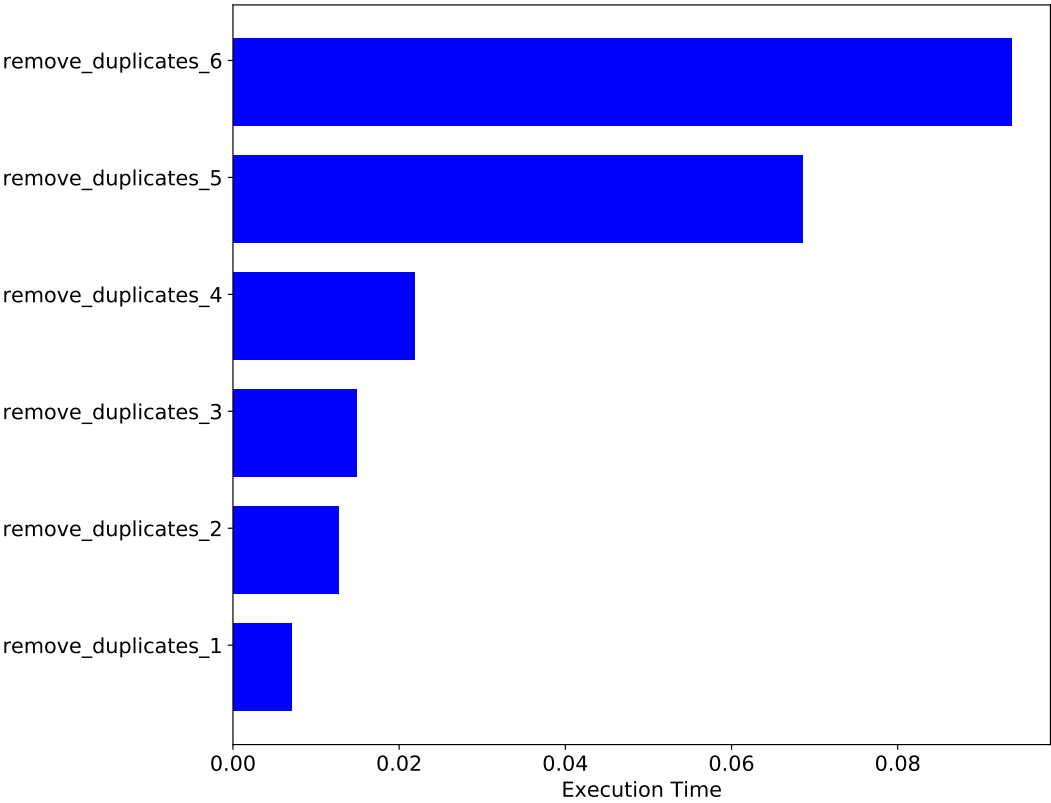
TIMES FOR THE FIRST 156000 NUMBERS



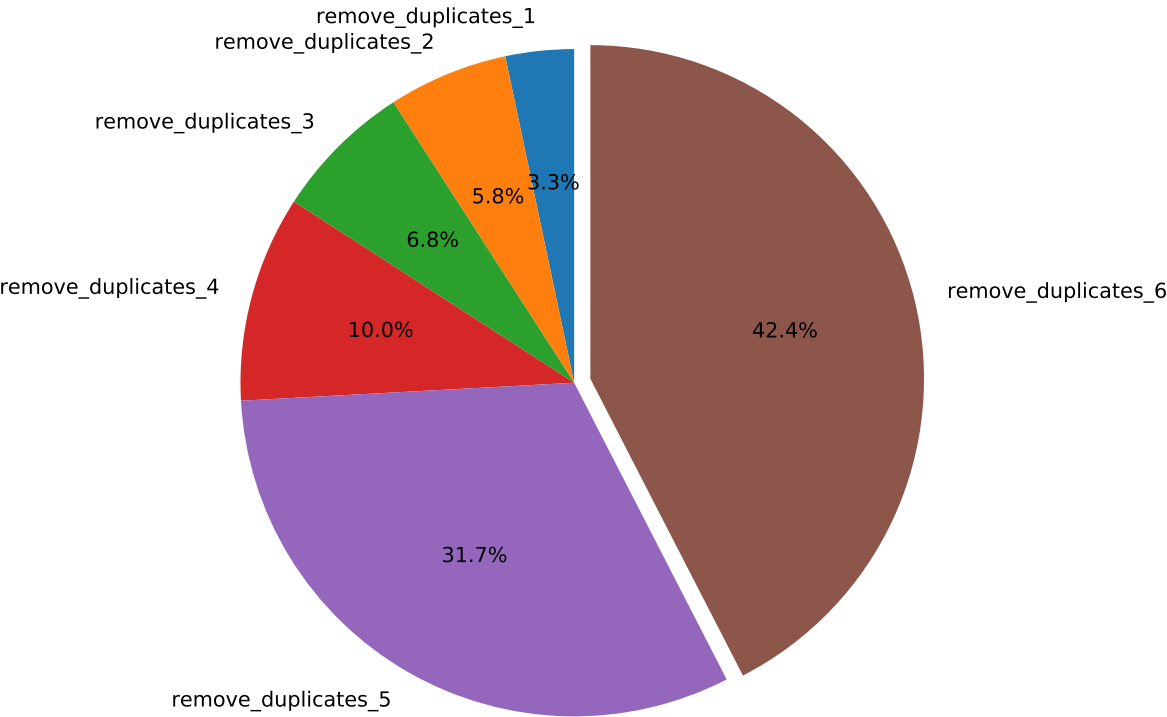
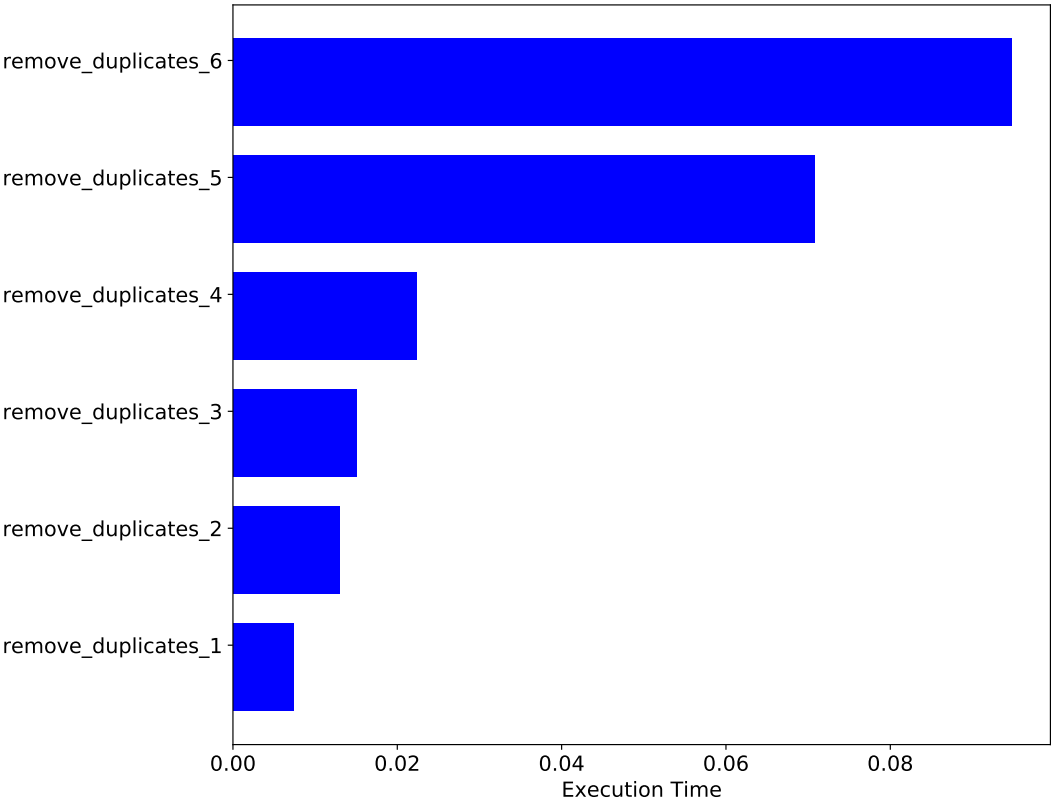
TIMES FOR THE FIRST 158000 NUMBERS



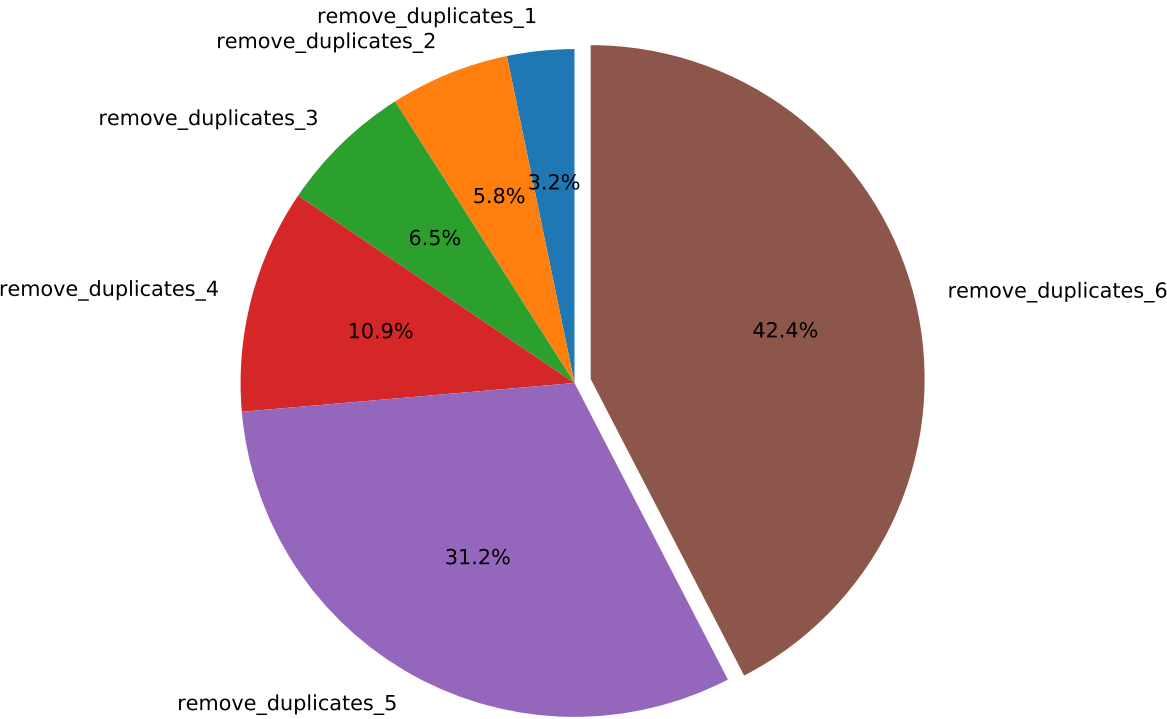
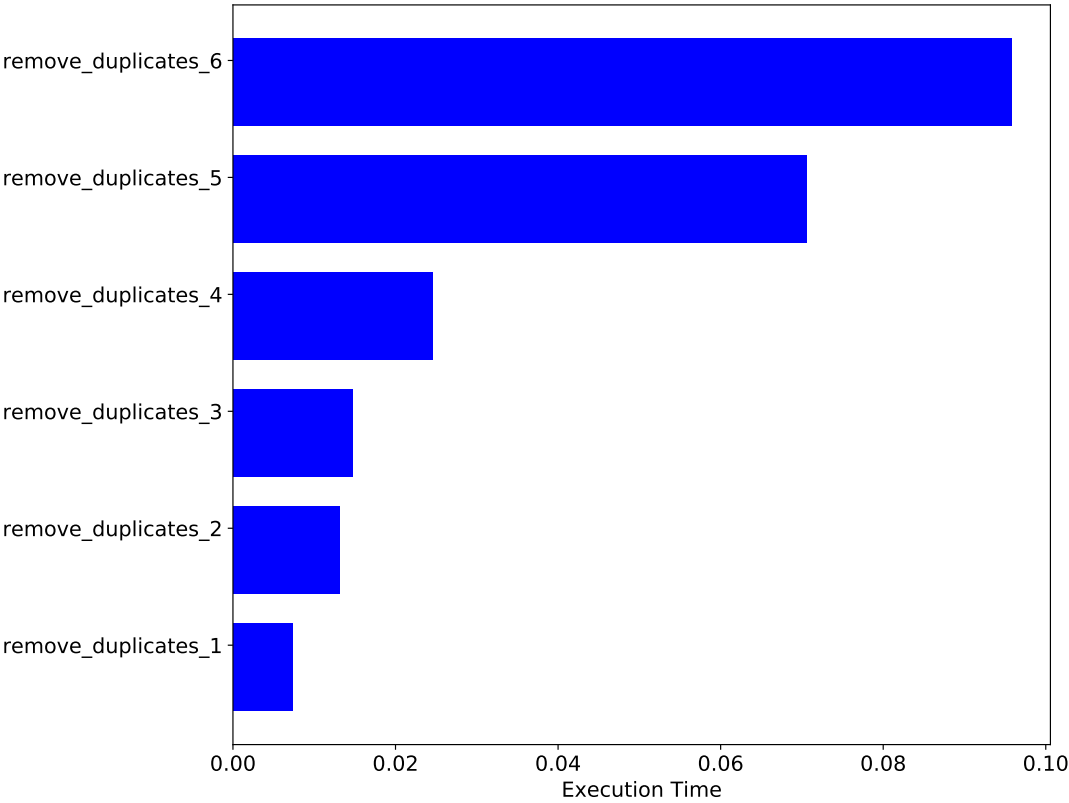
TIMES FOR THE FIRST 160000 NUMBERS



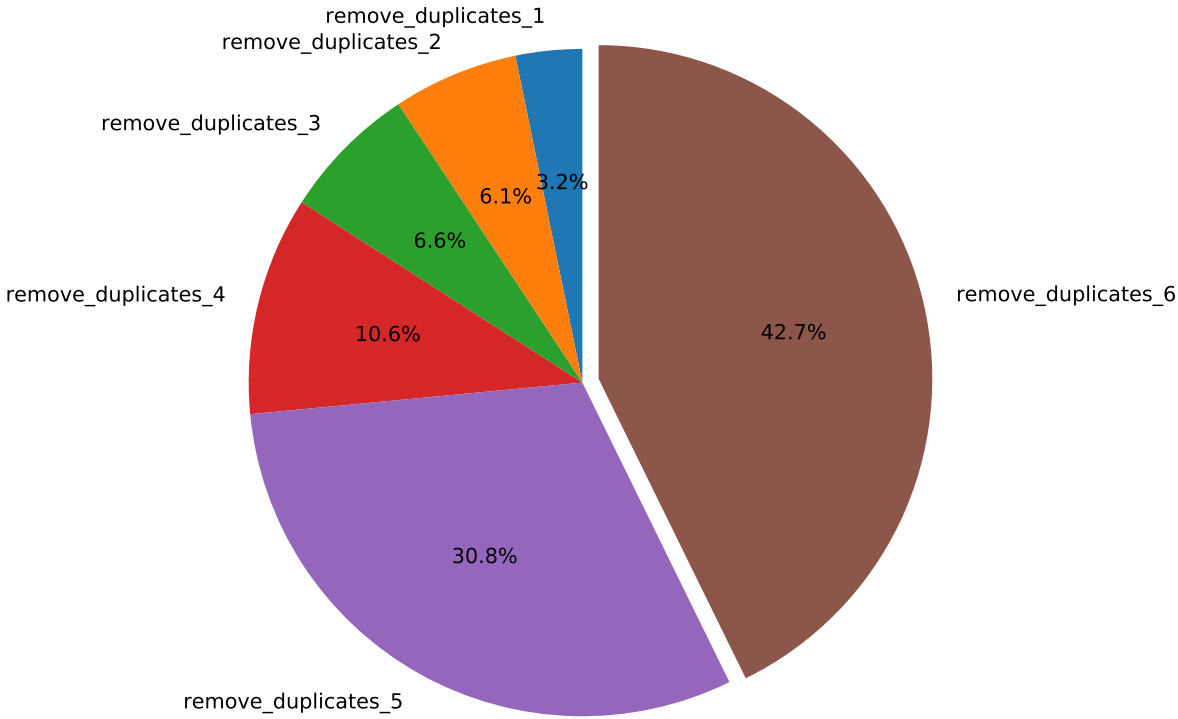
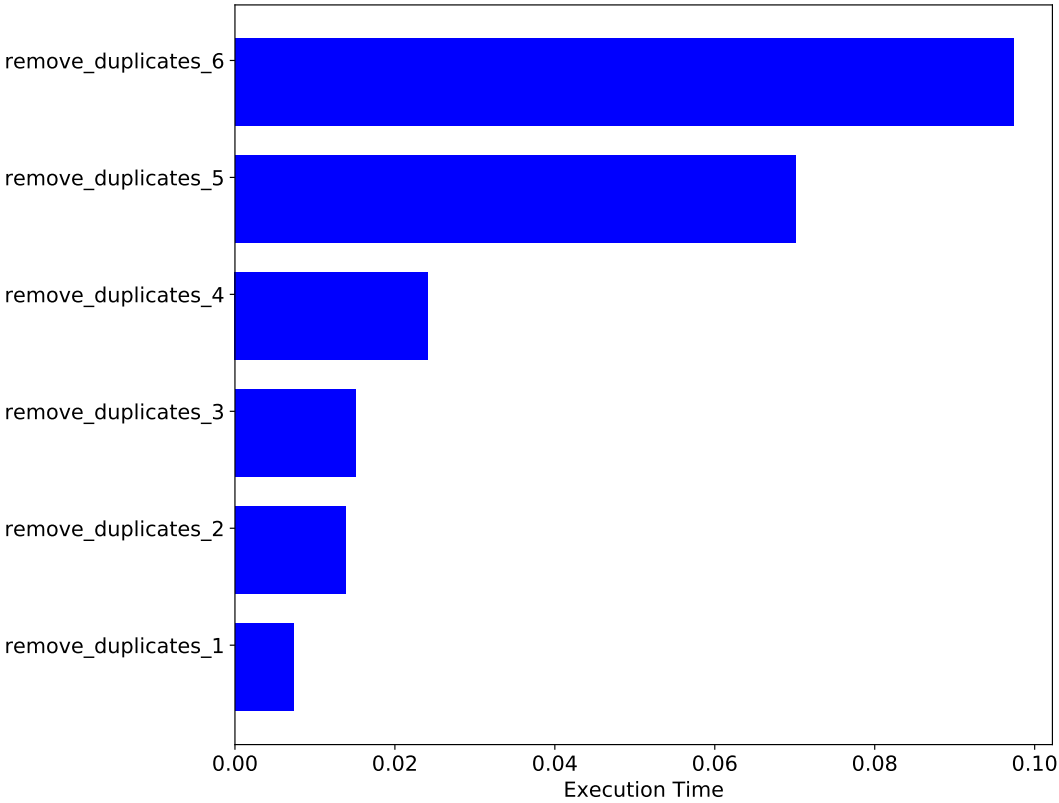
TIMES FOR THE FIRST 162000 NUMBERS



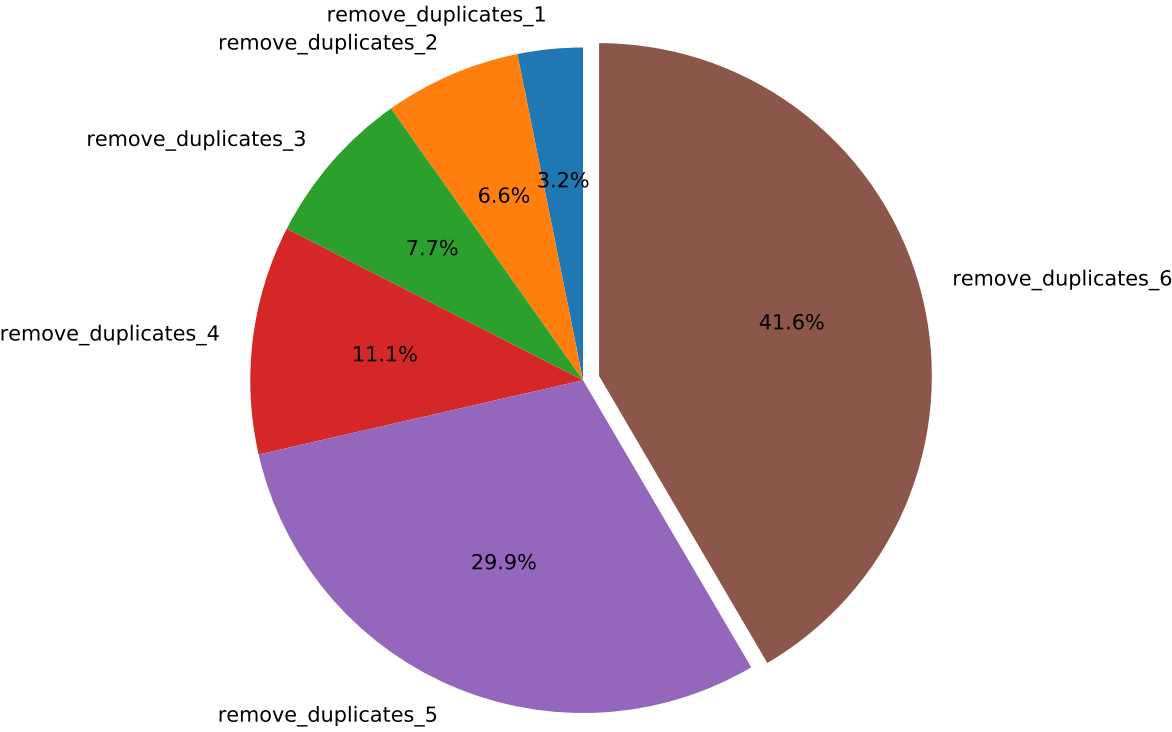
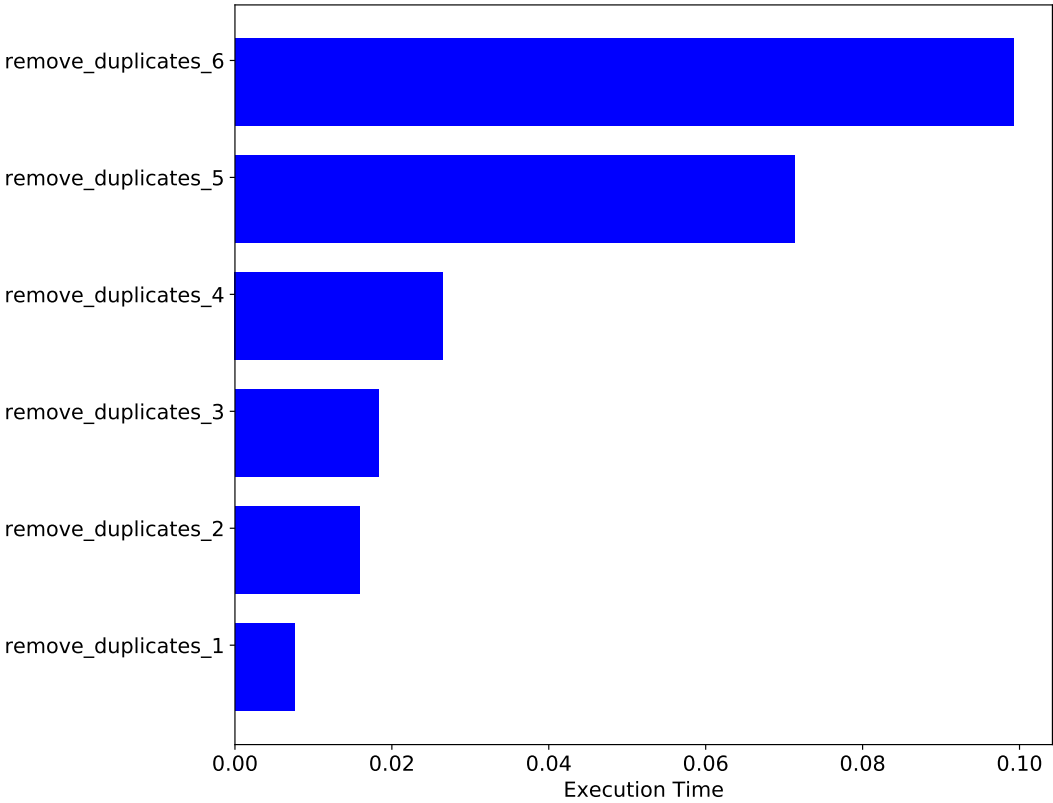
TIMES FOR THE FIRST 164000 NUMBERS



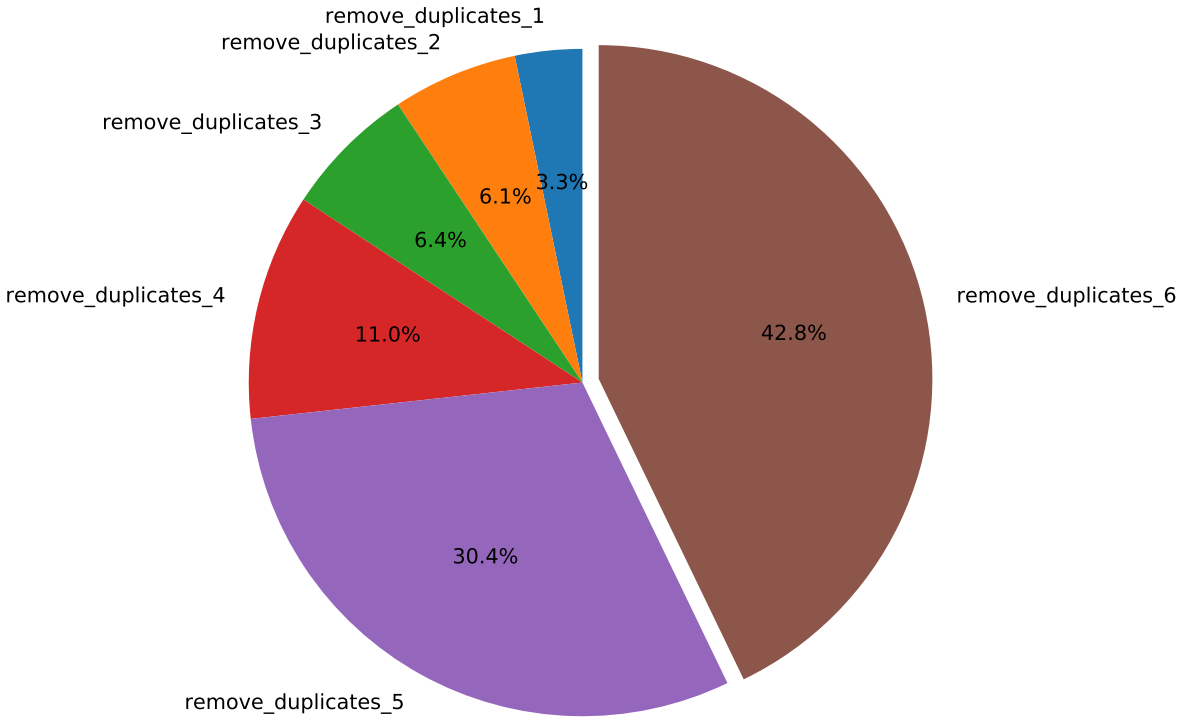
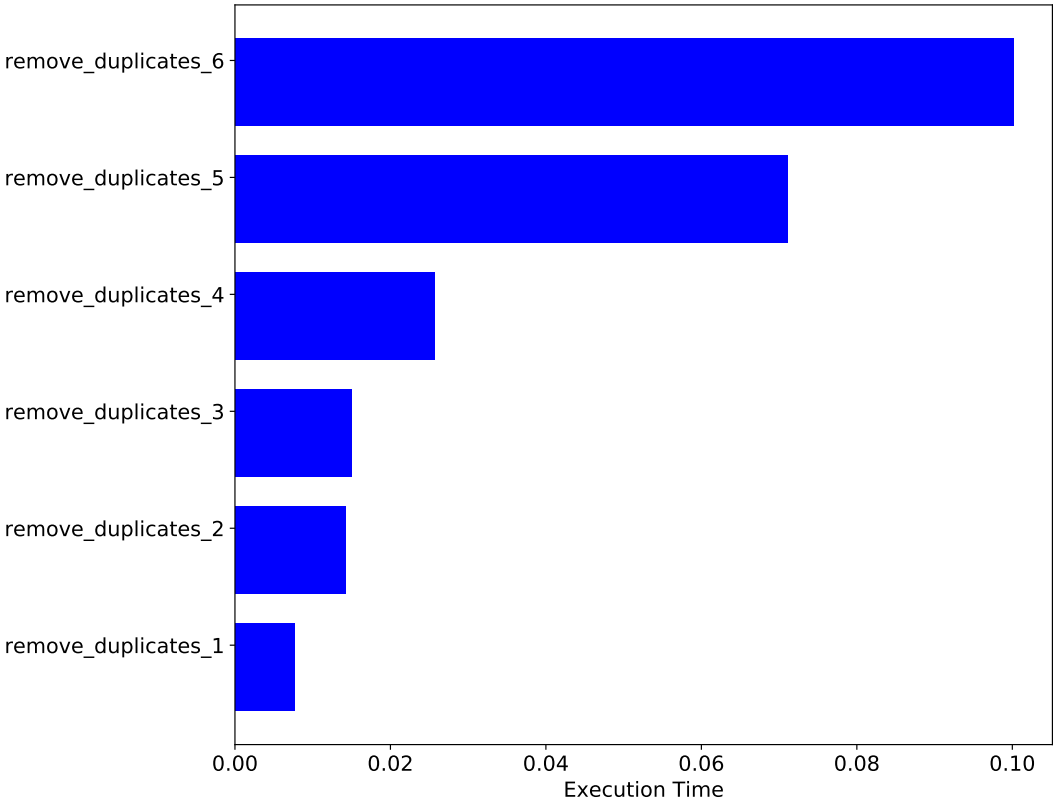
TIMES FOR THE FIRST 166000 NUMBERS



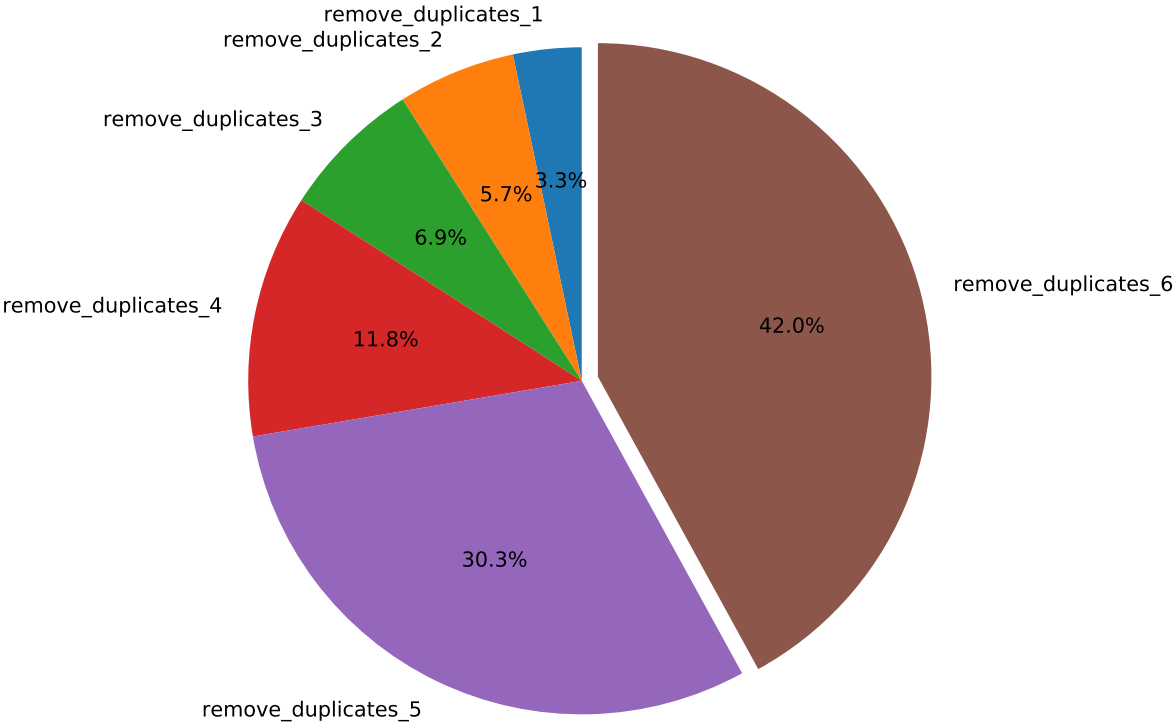
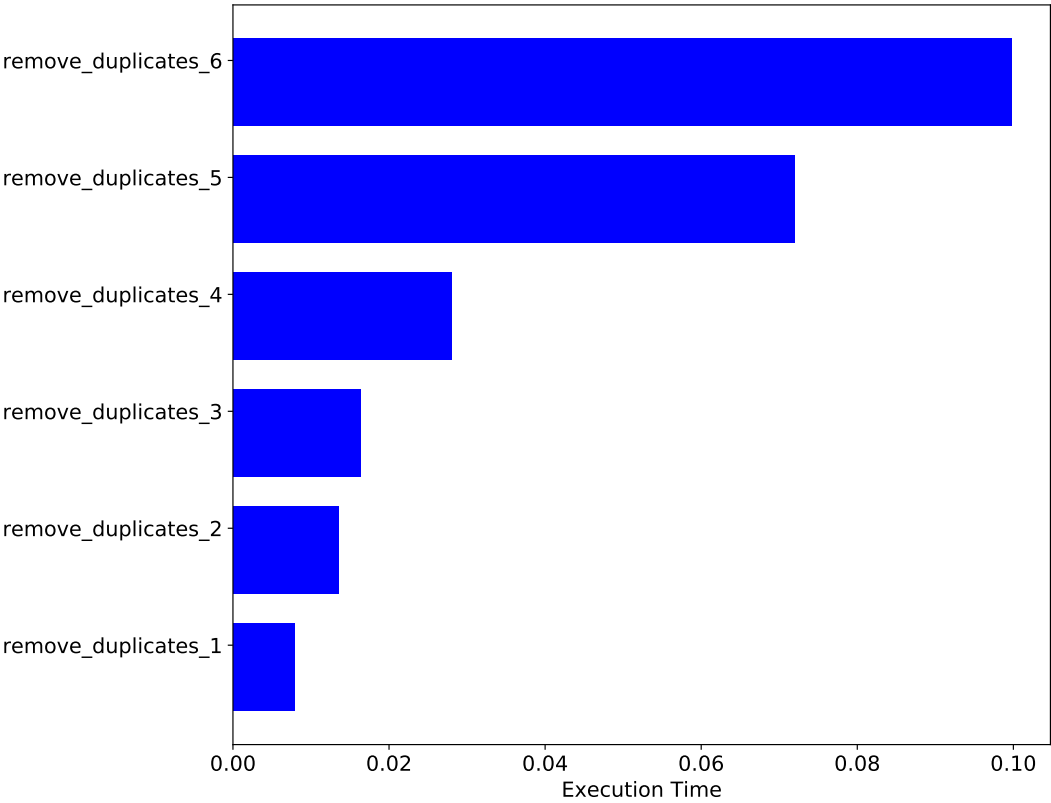
TIMES FOR THE FIRST 168000 NUMBERS



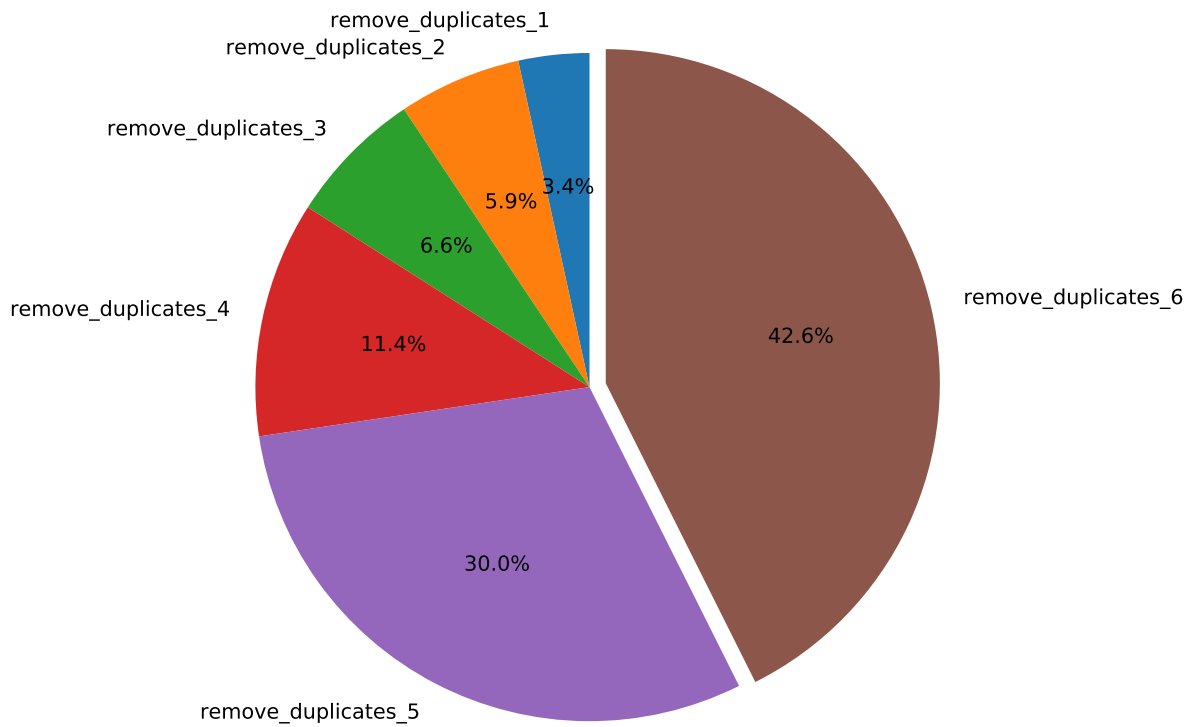
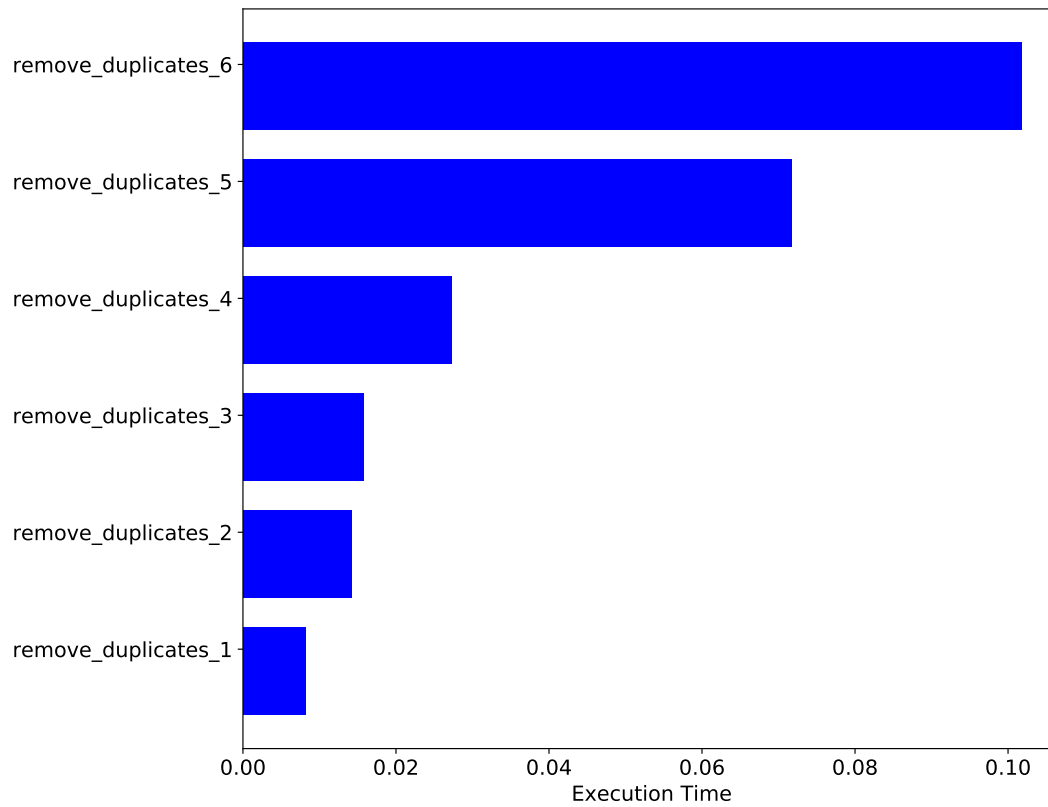
TIMES FOR THE FIRST 170000 NUMBERS



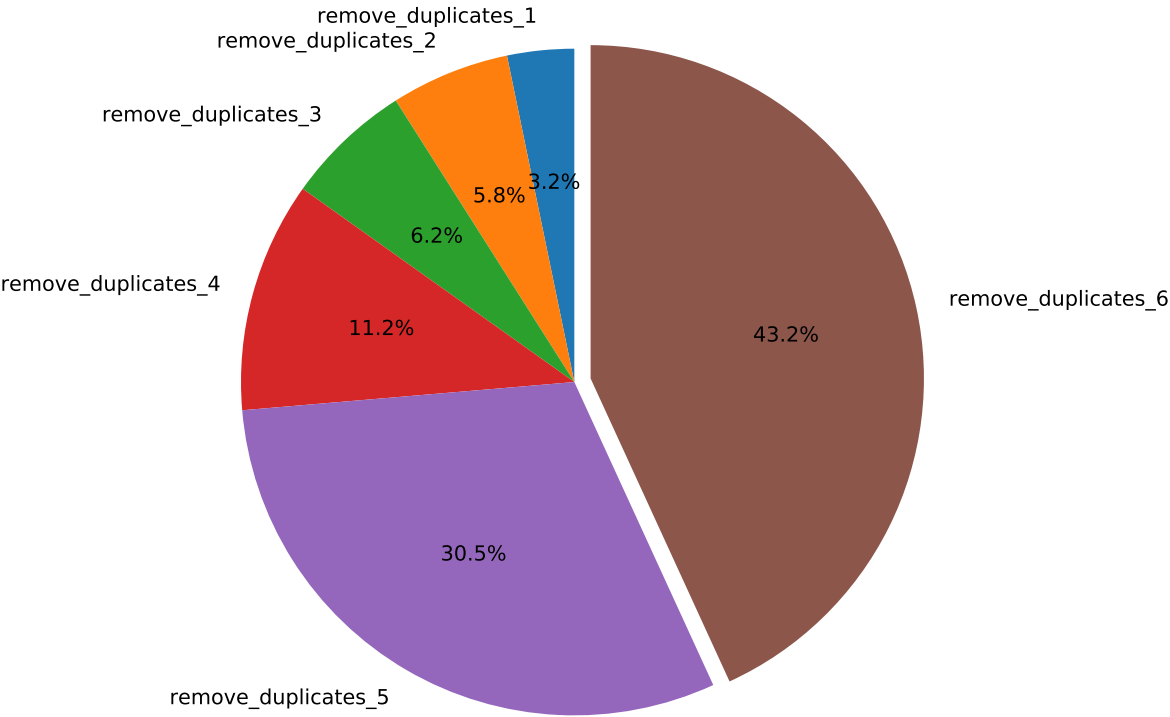
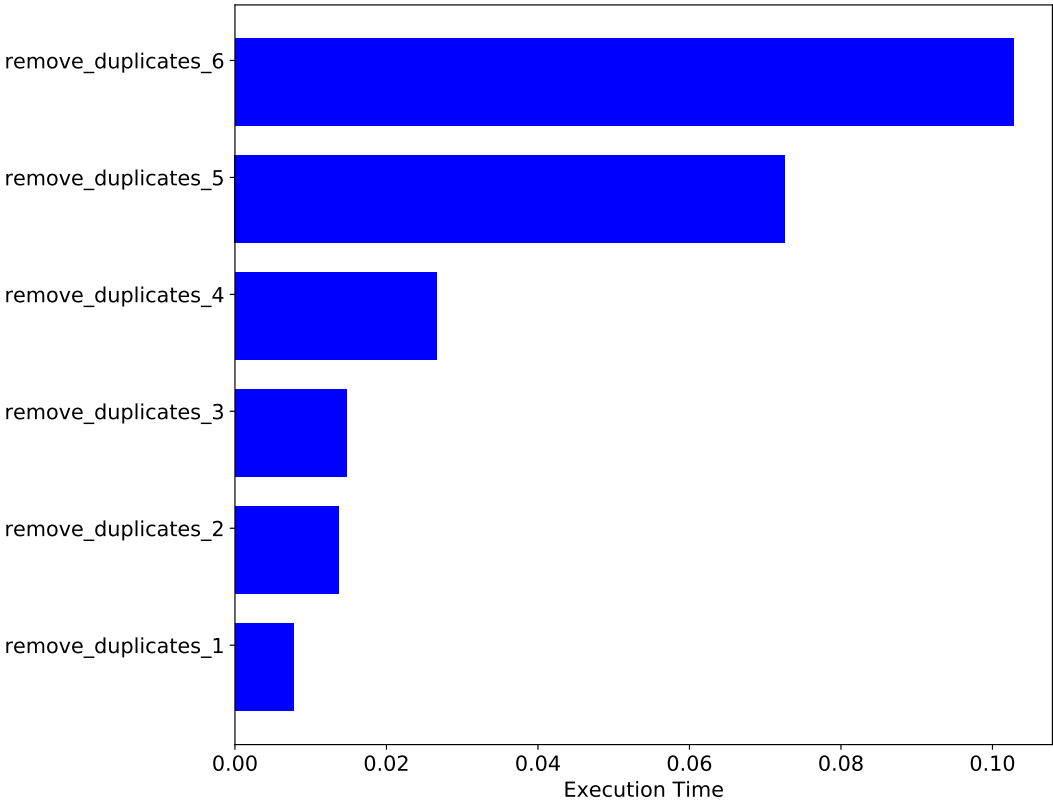
TIMES FOR THE FIRST 172000 NUMBERS



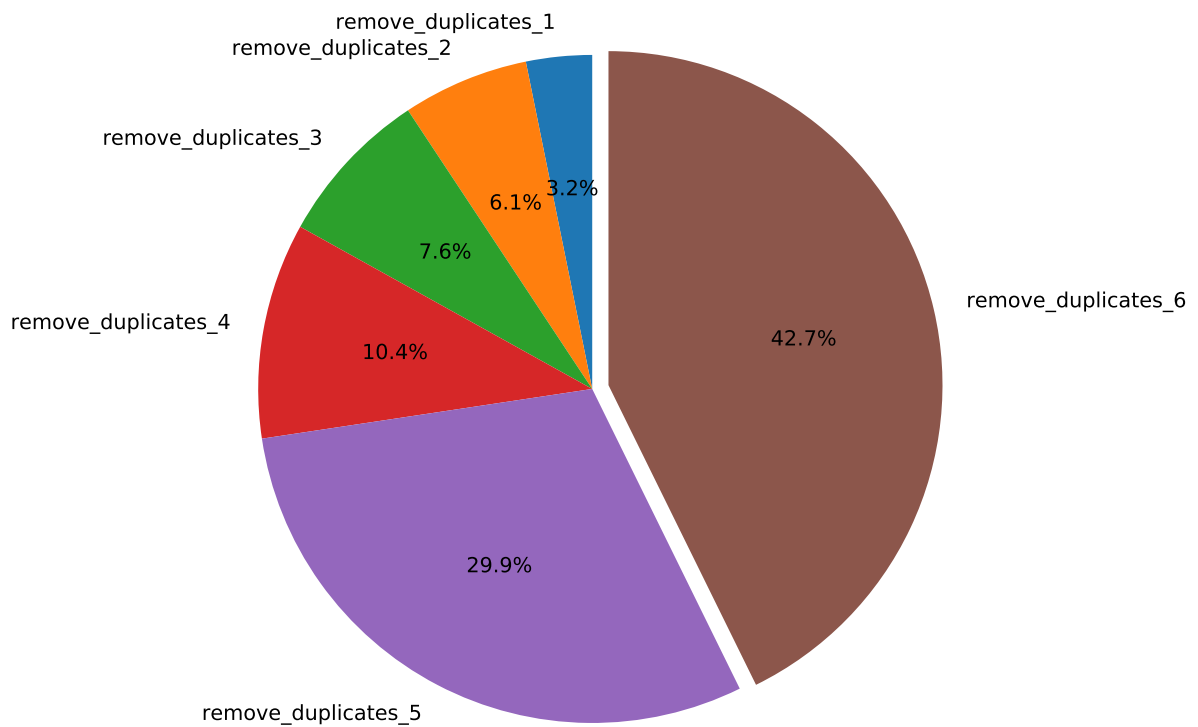
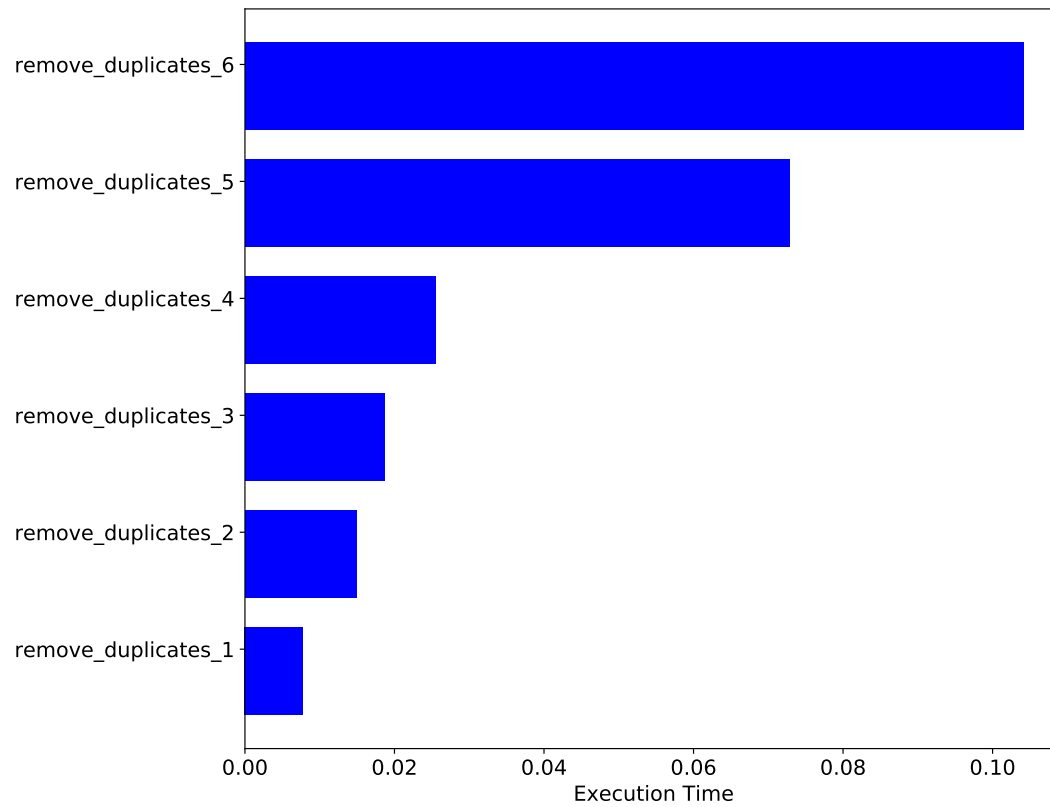
TIMES FOR THE FIRST 174000 NUMBERS



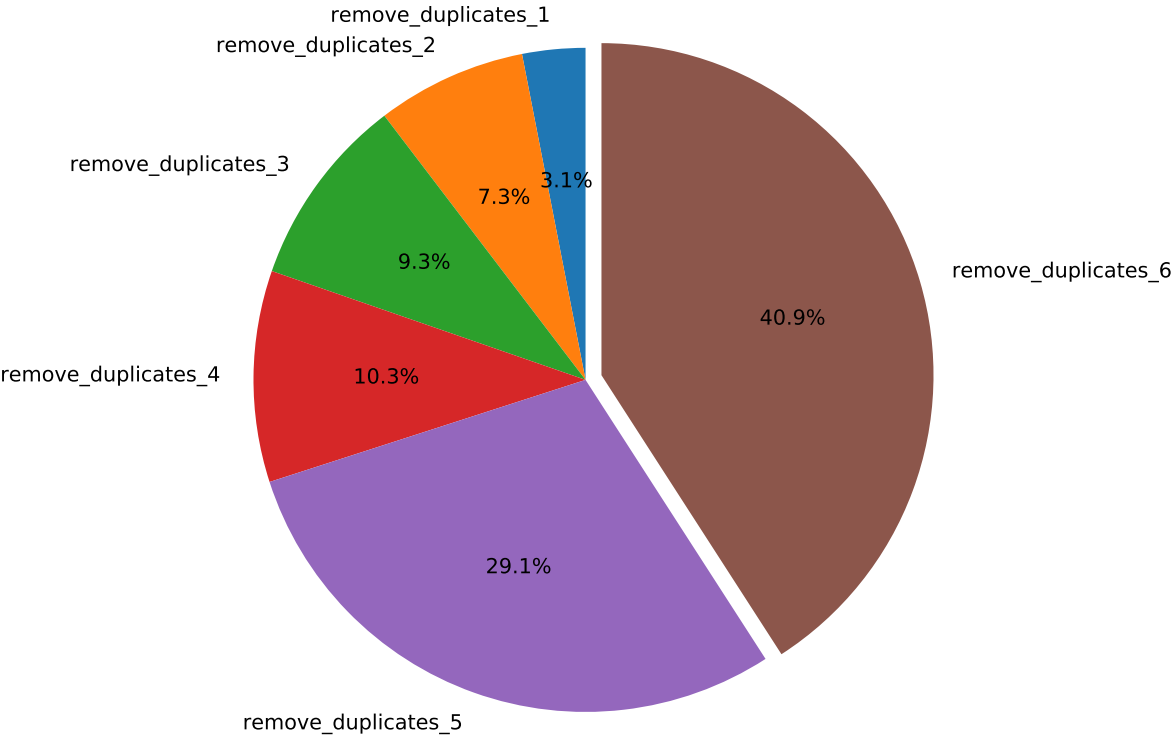
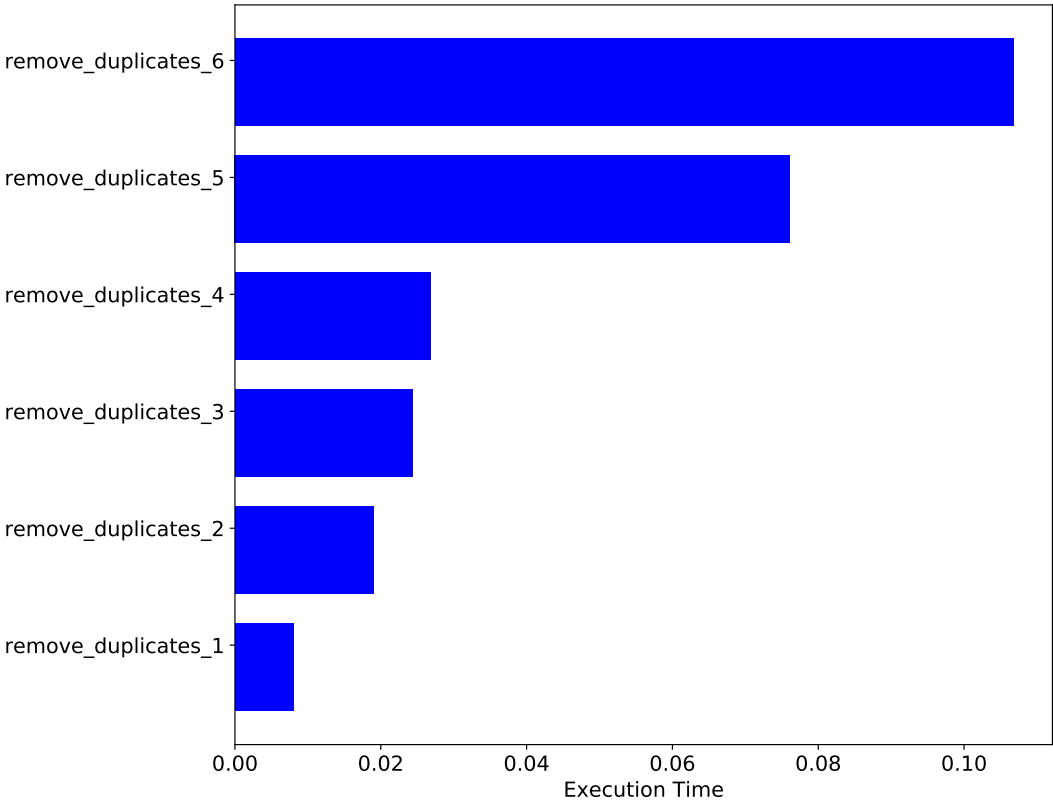
TIMES FOR THE FIRST 176000 NUMBERS



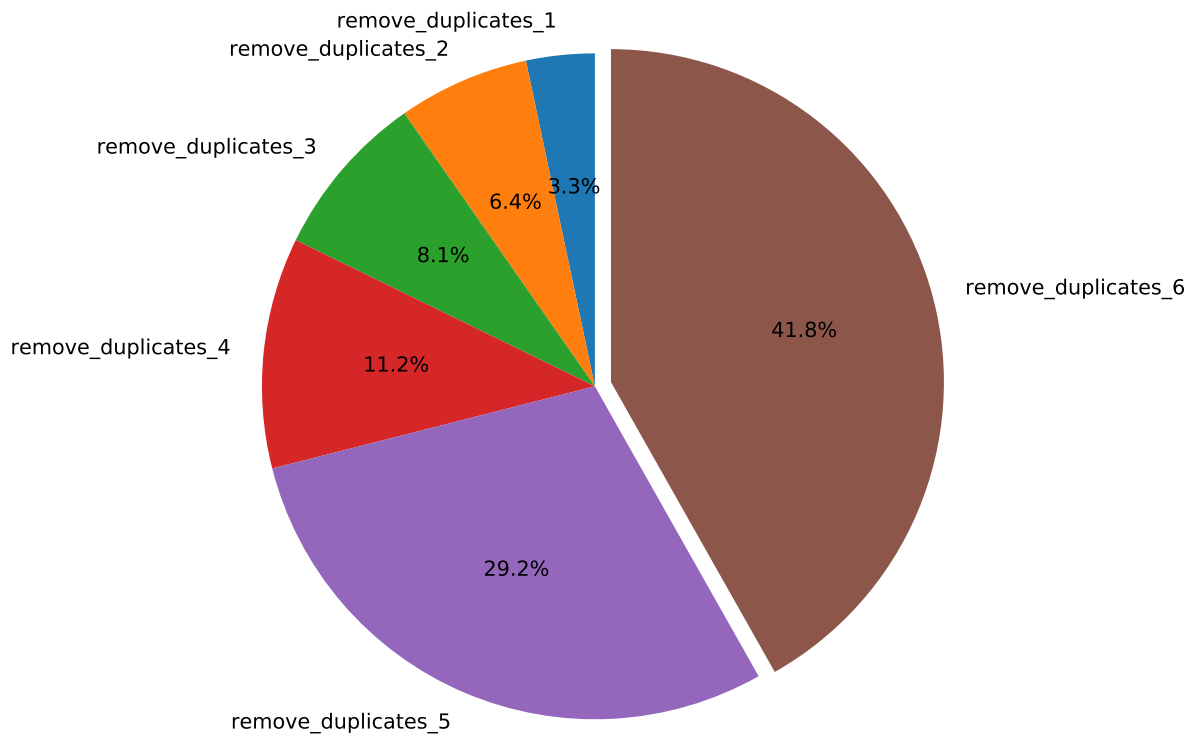
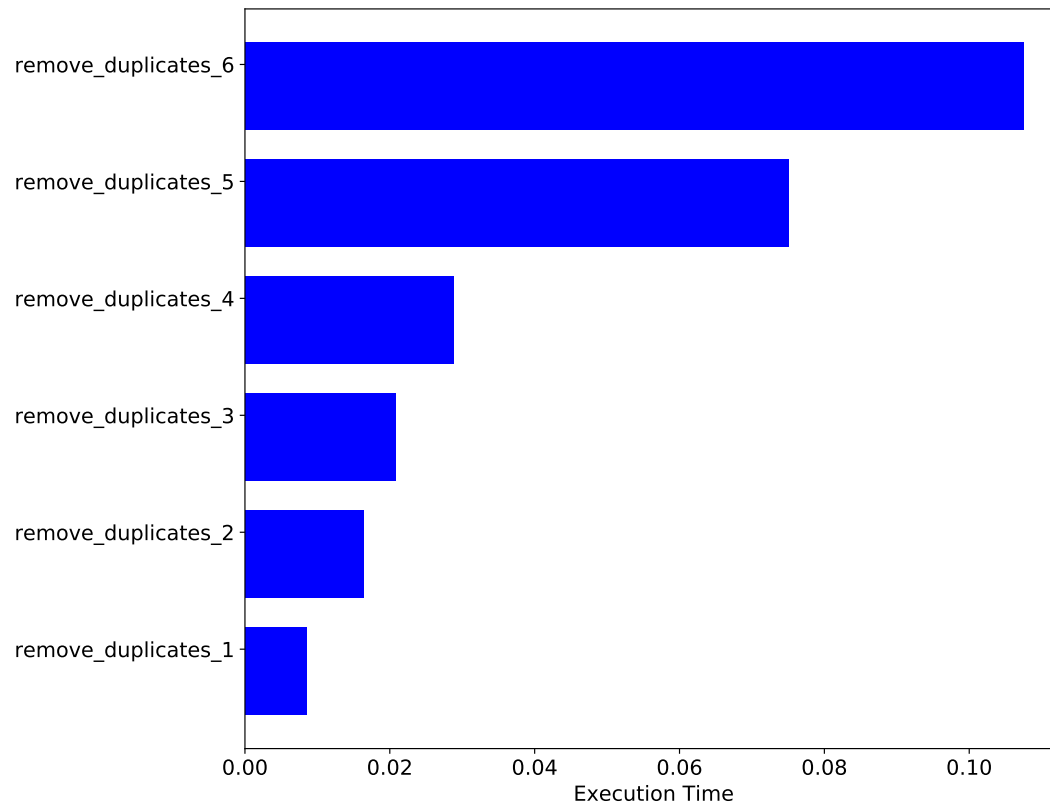
TIMES FOR THE FIRST 178000 NUMBERS



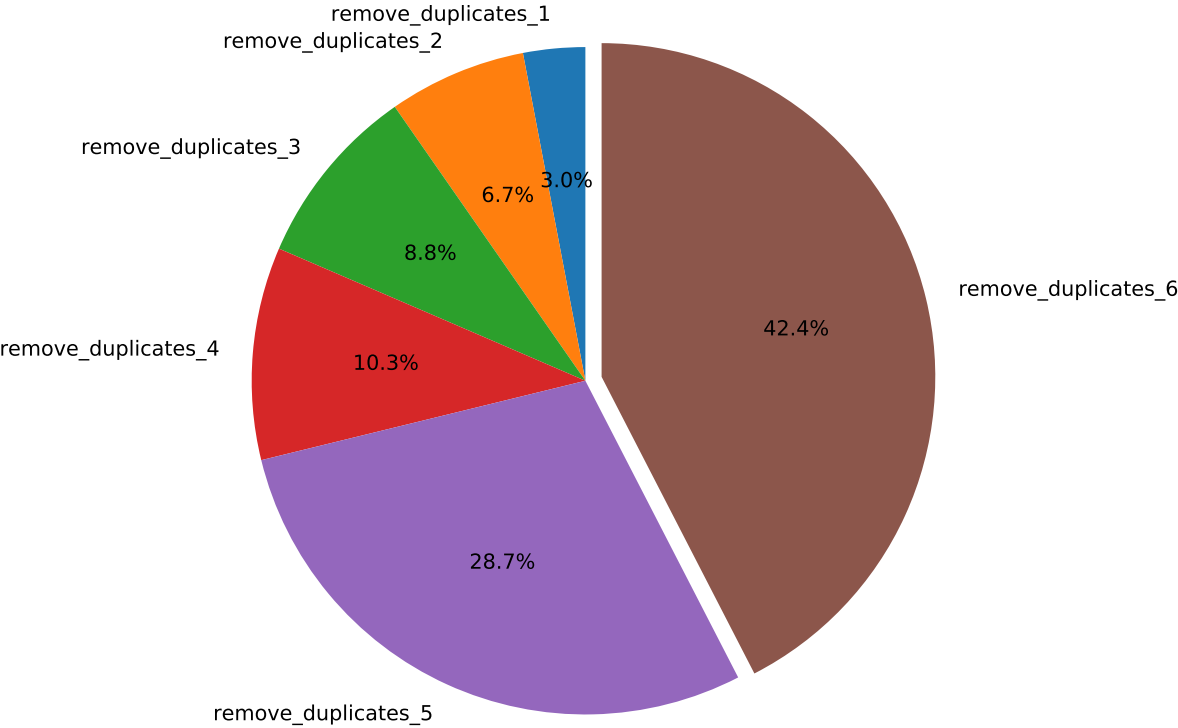
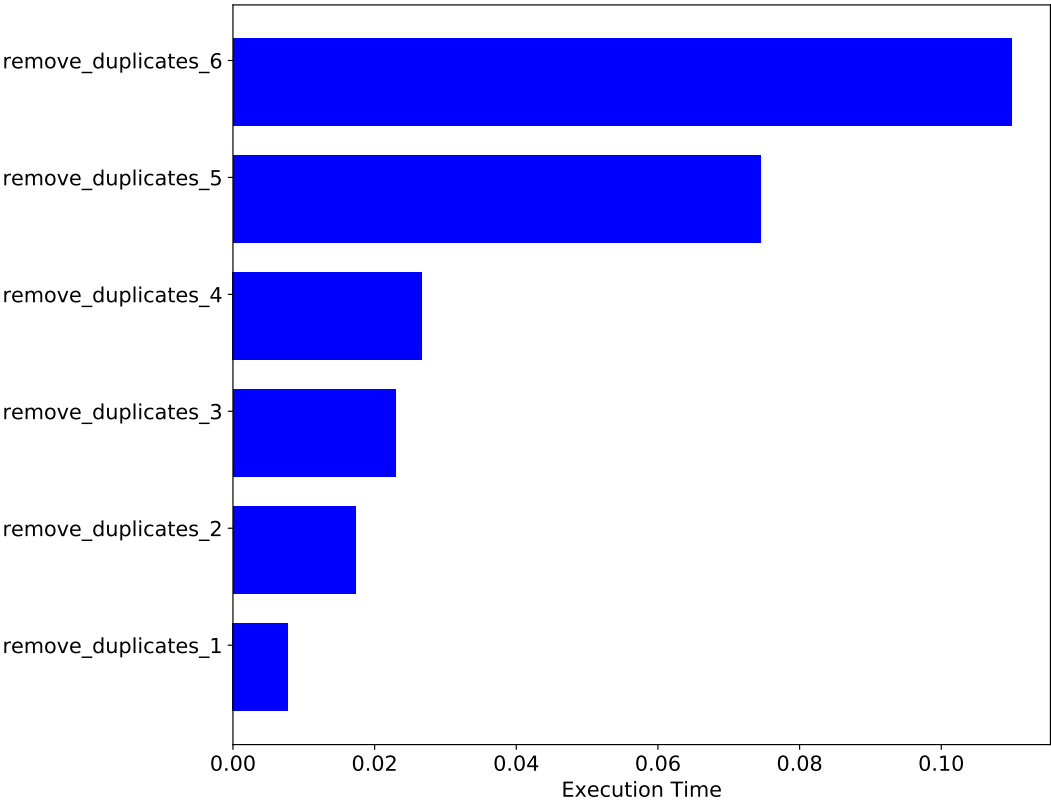
TIMES FOR THE FIRST 180000 NUMBERS



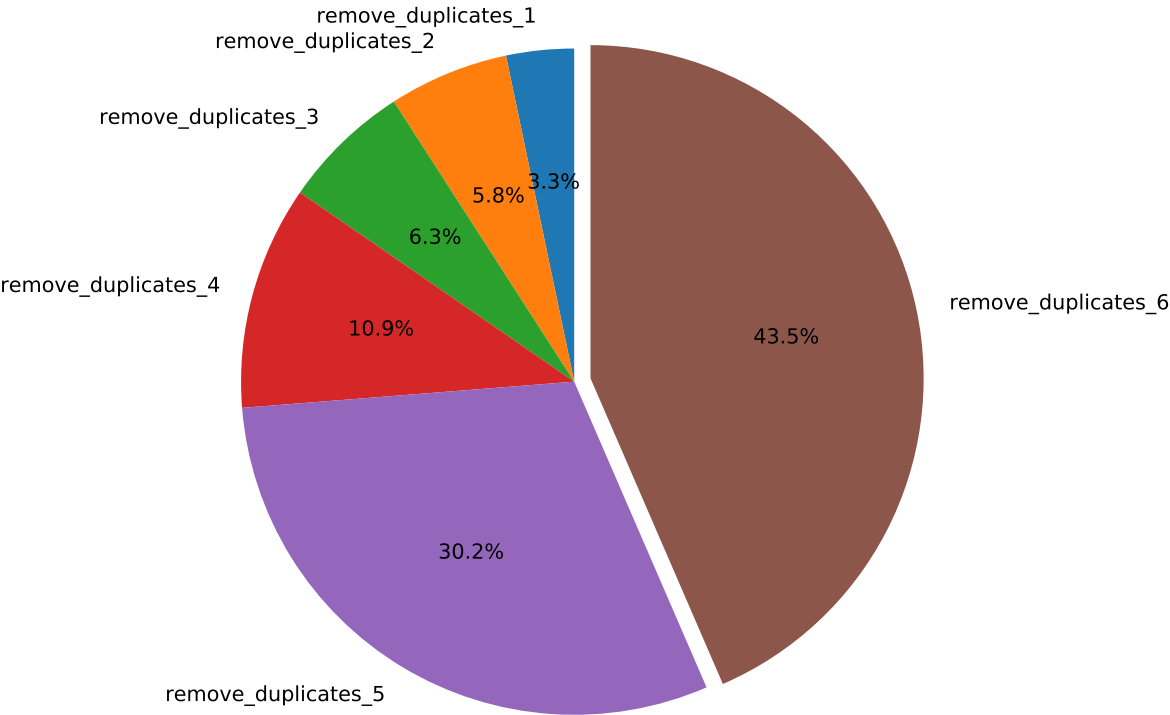
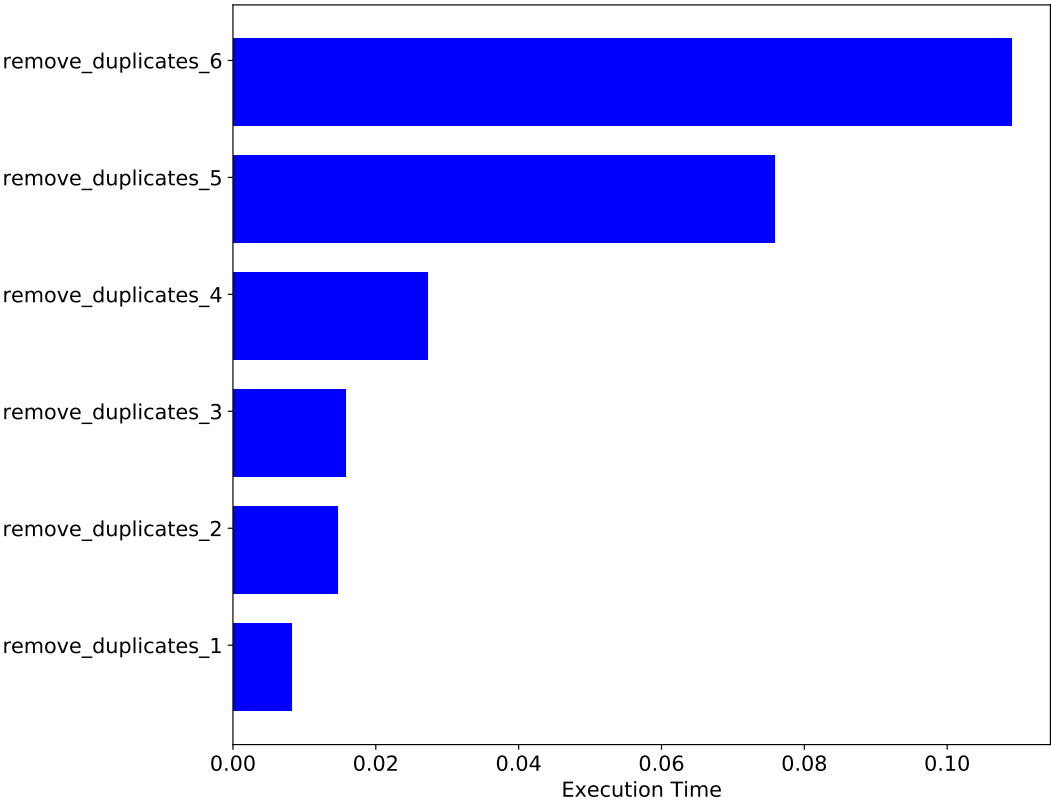
TIMES FOR THE FIRST 182000 NUMBERS



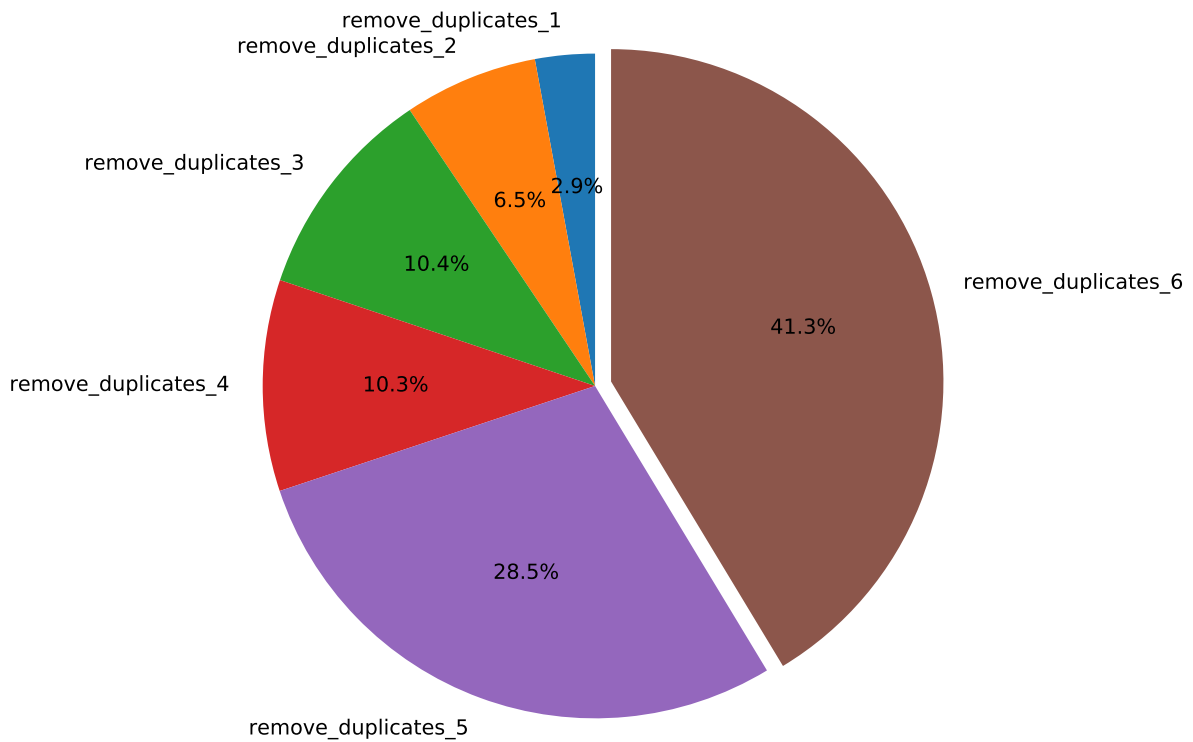
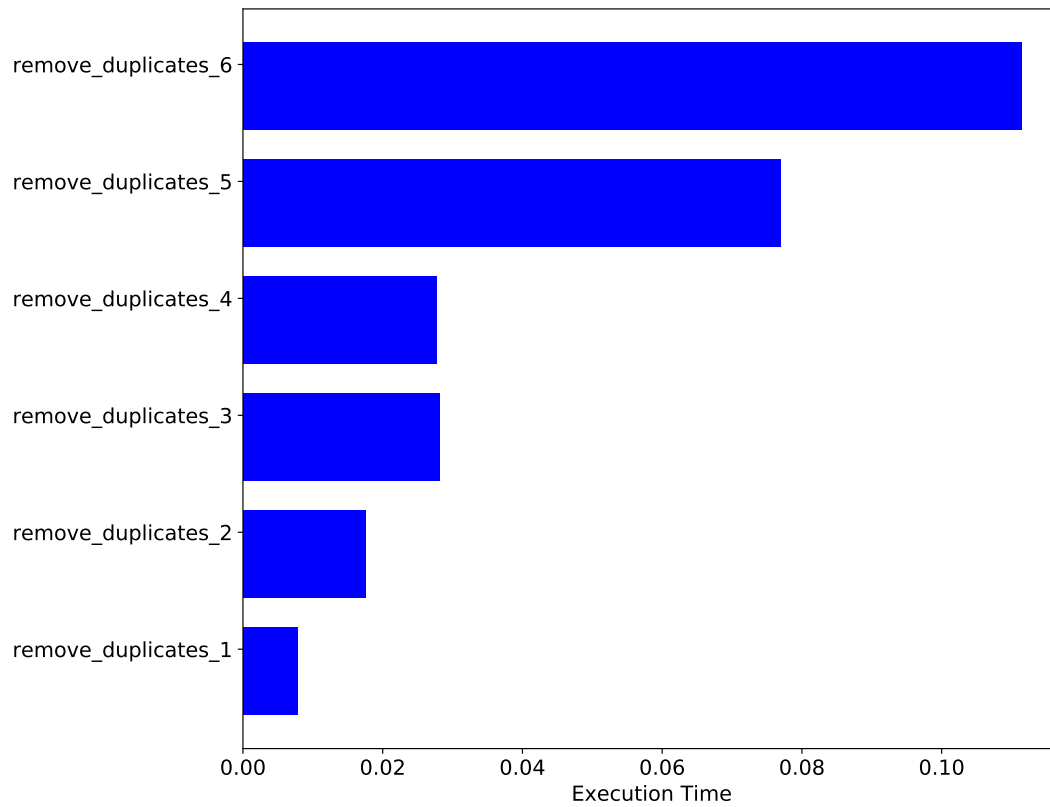
TIMES FOR THE FIRST 184000 NUMBERS



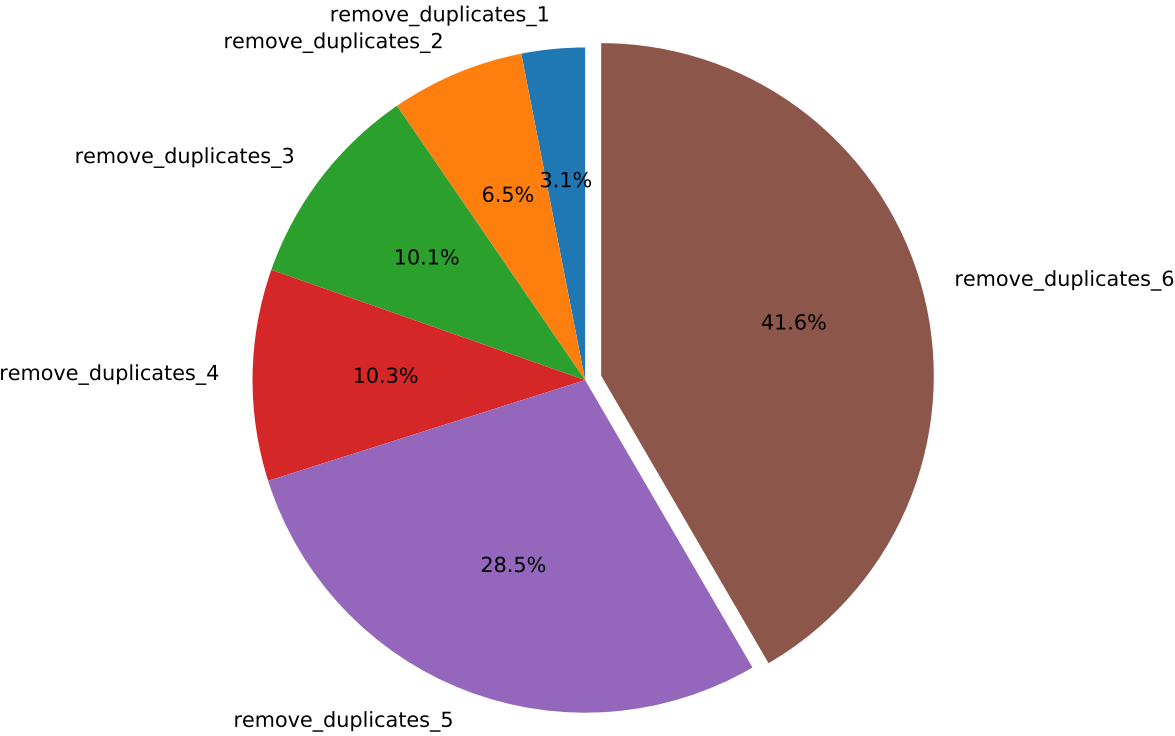
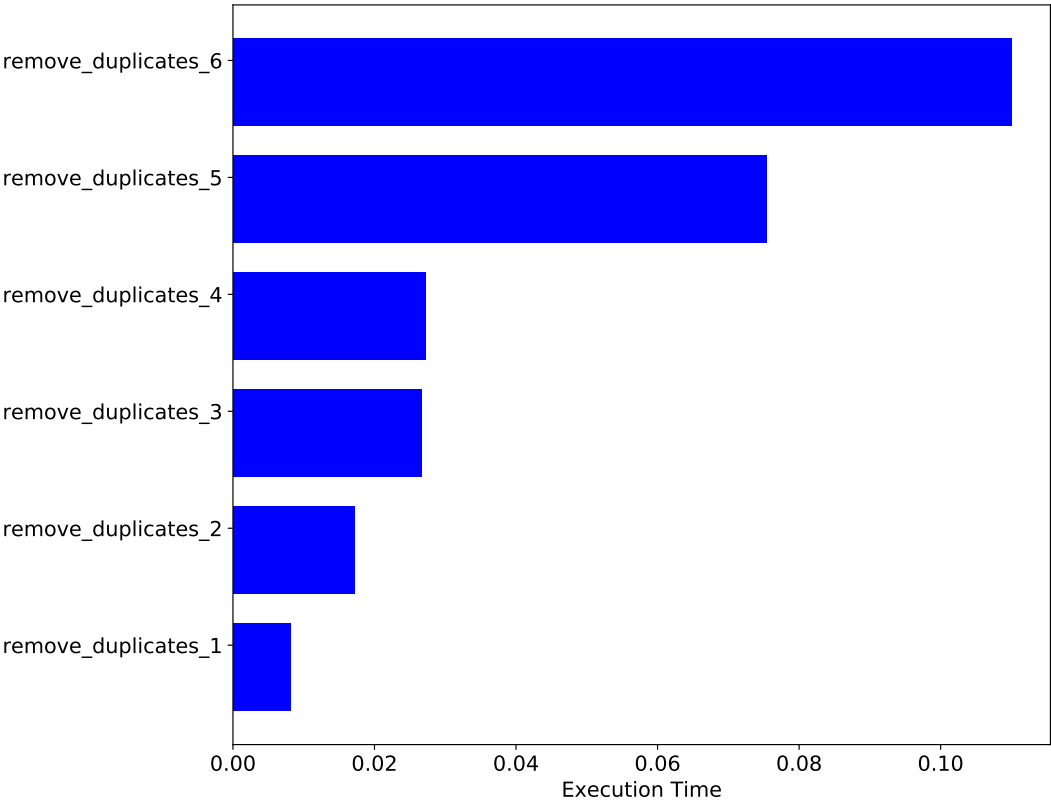
TIMES FOR THE FIRST 186000 NUMBERS



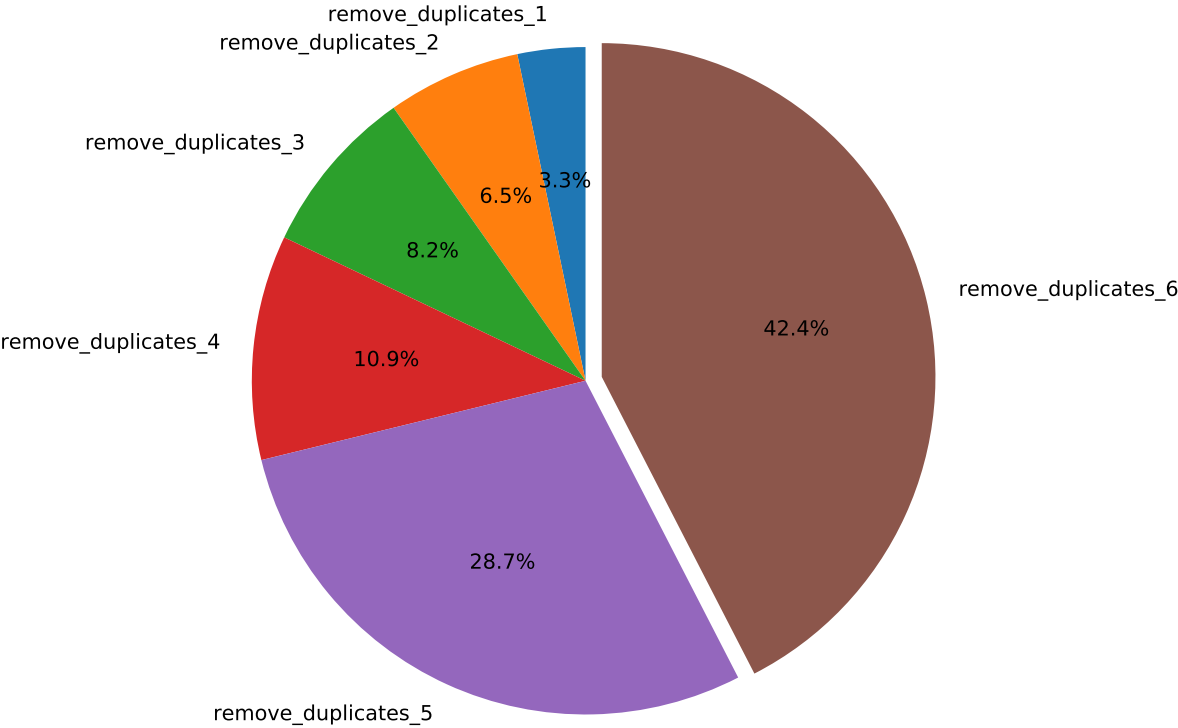
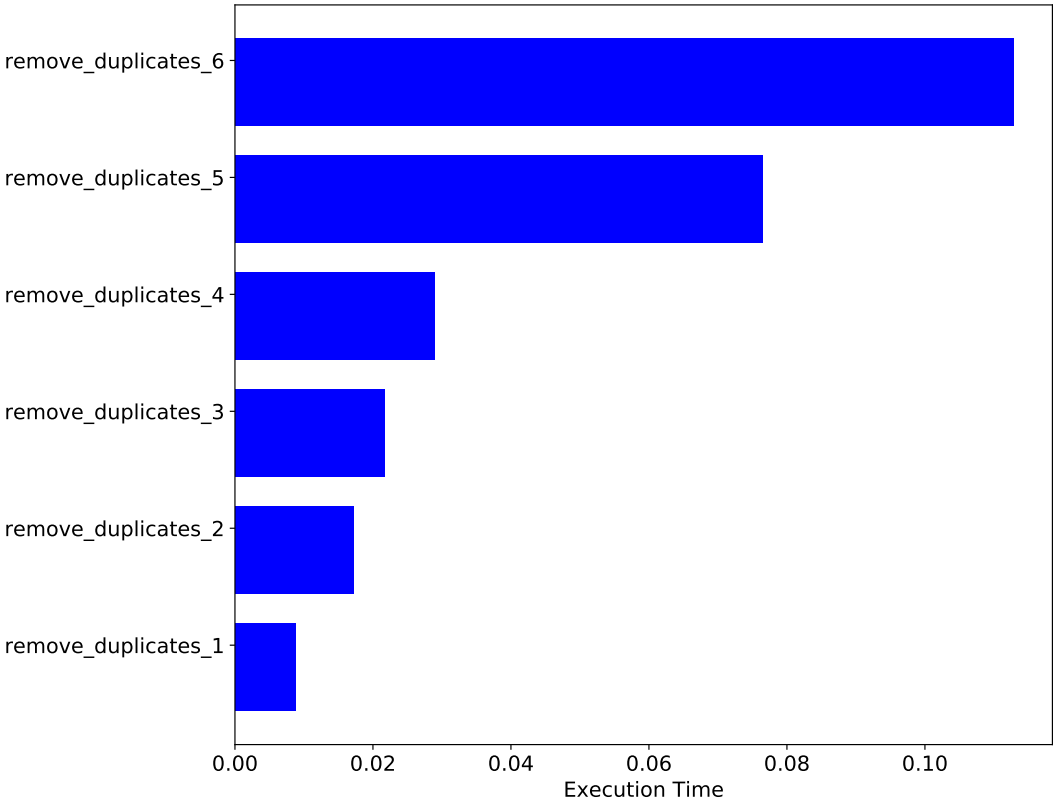
TIMES FOR THE FIRST 188000 NUMBERS



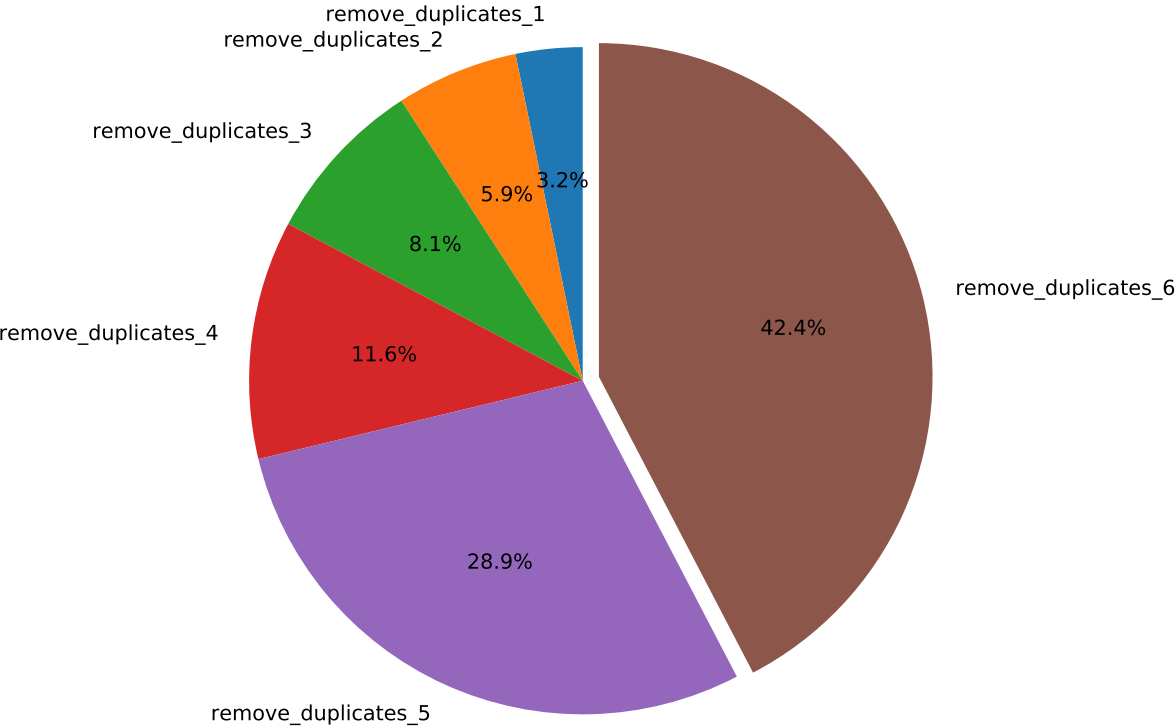
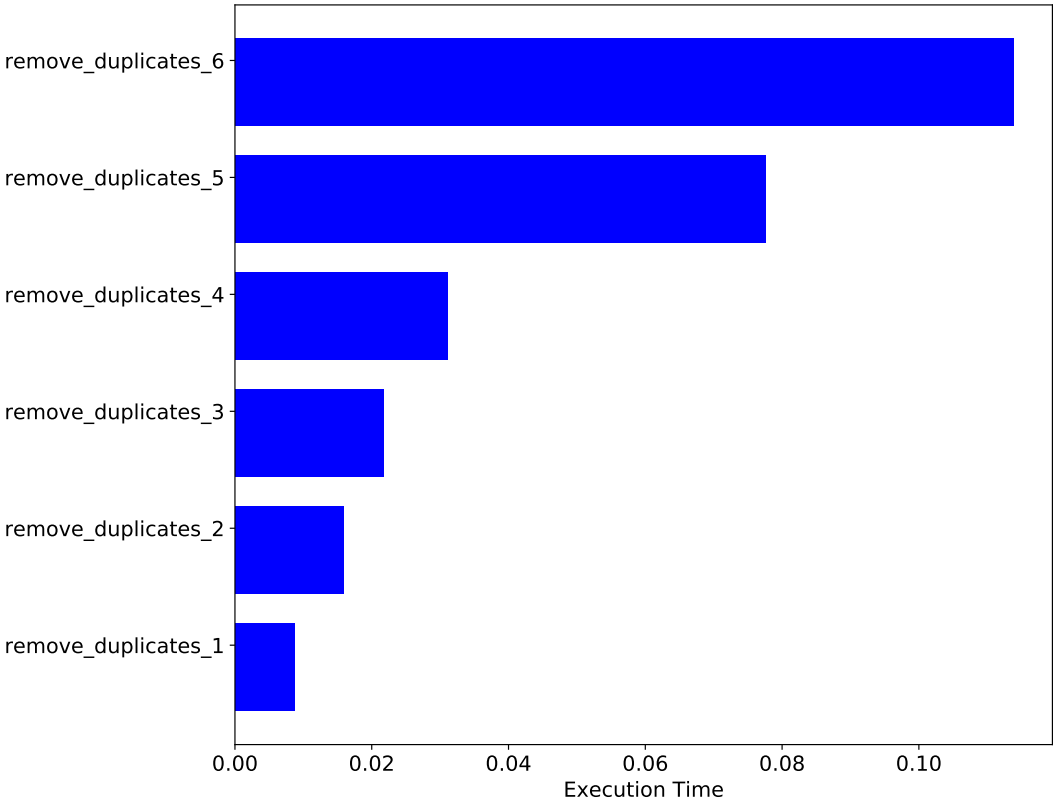
TIMES FOR THE FIRST 190000 NUMBERS



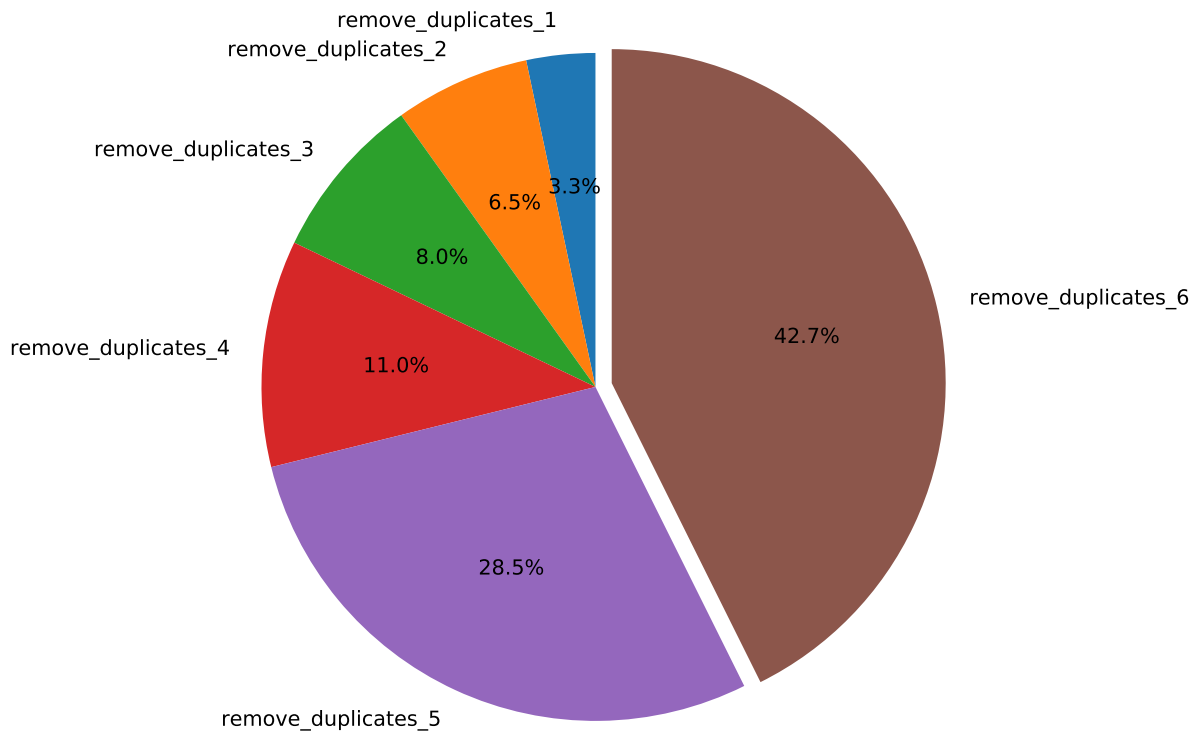
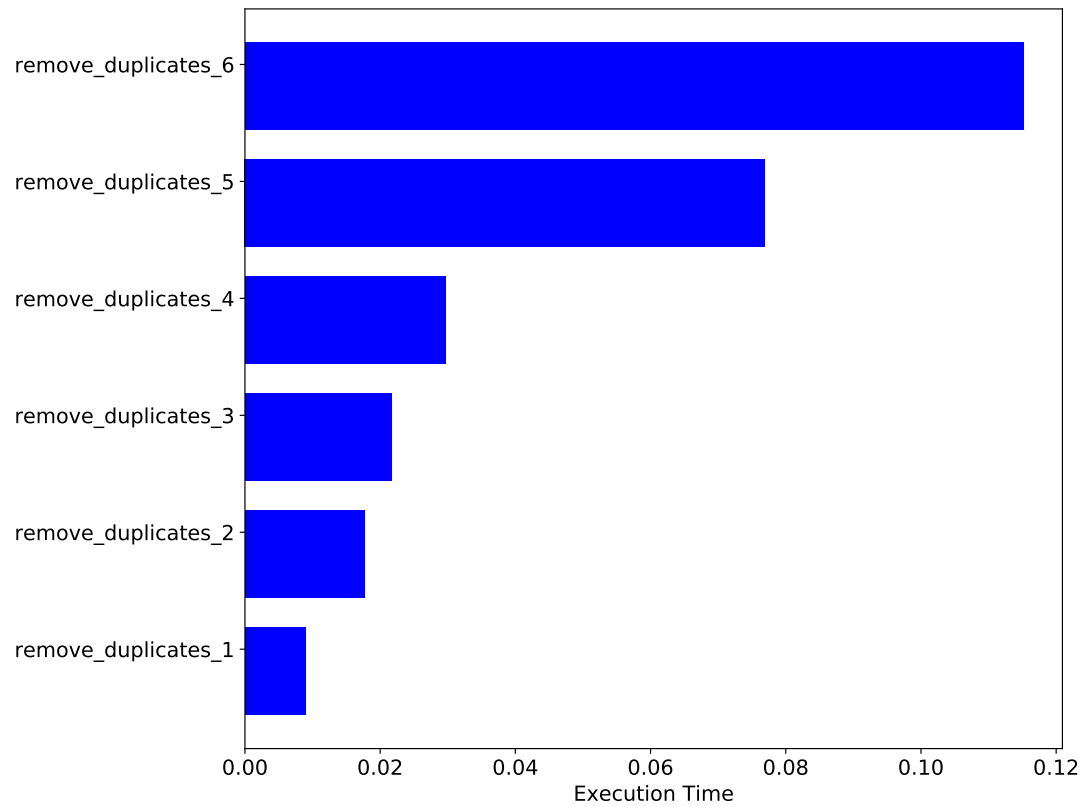
TIMES FOR THE FIRST 192000 NUMBERS



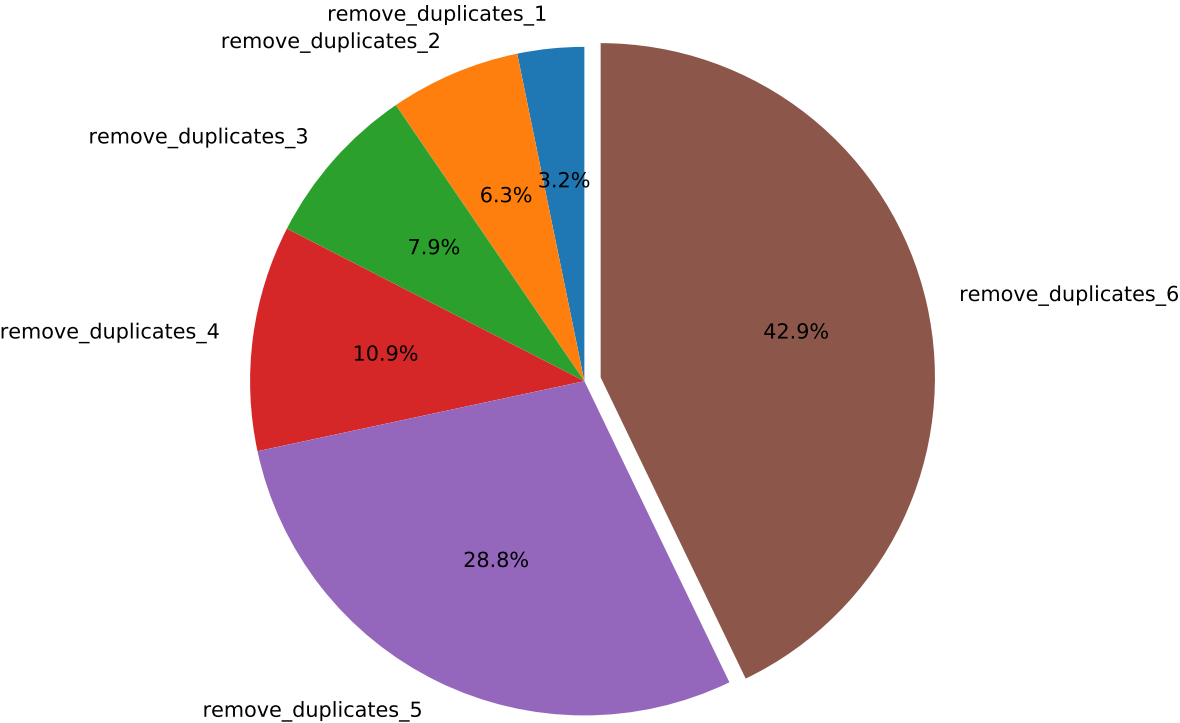
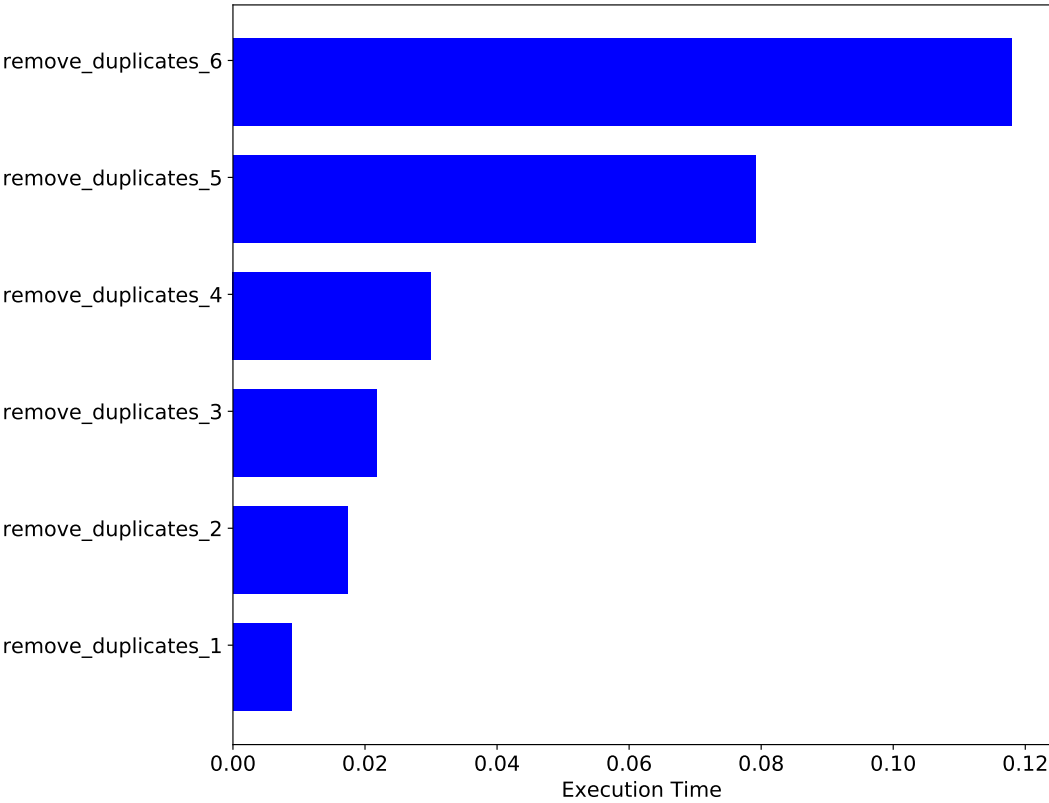
TIMES FOR THE FIRST 194000 NUMBERS



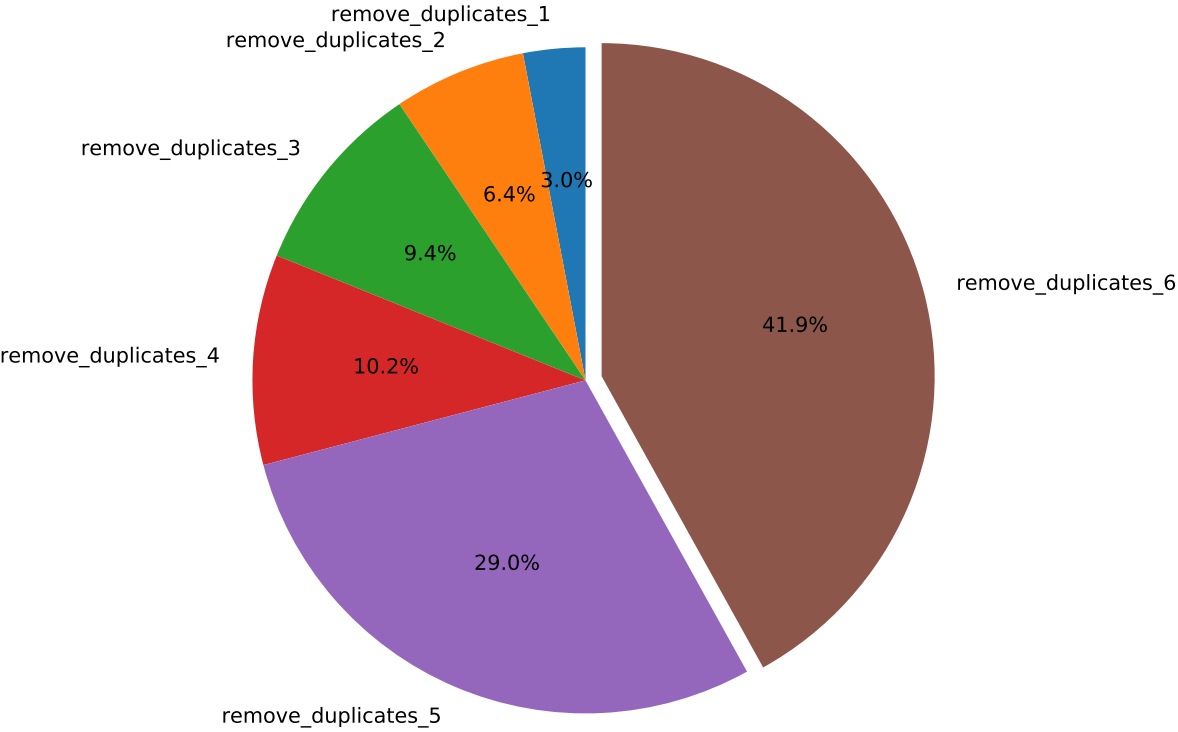
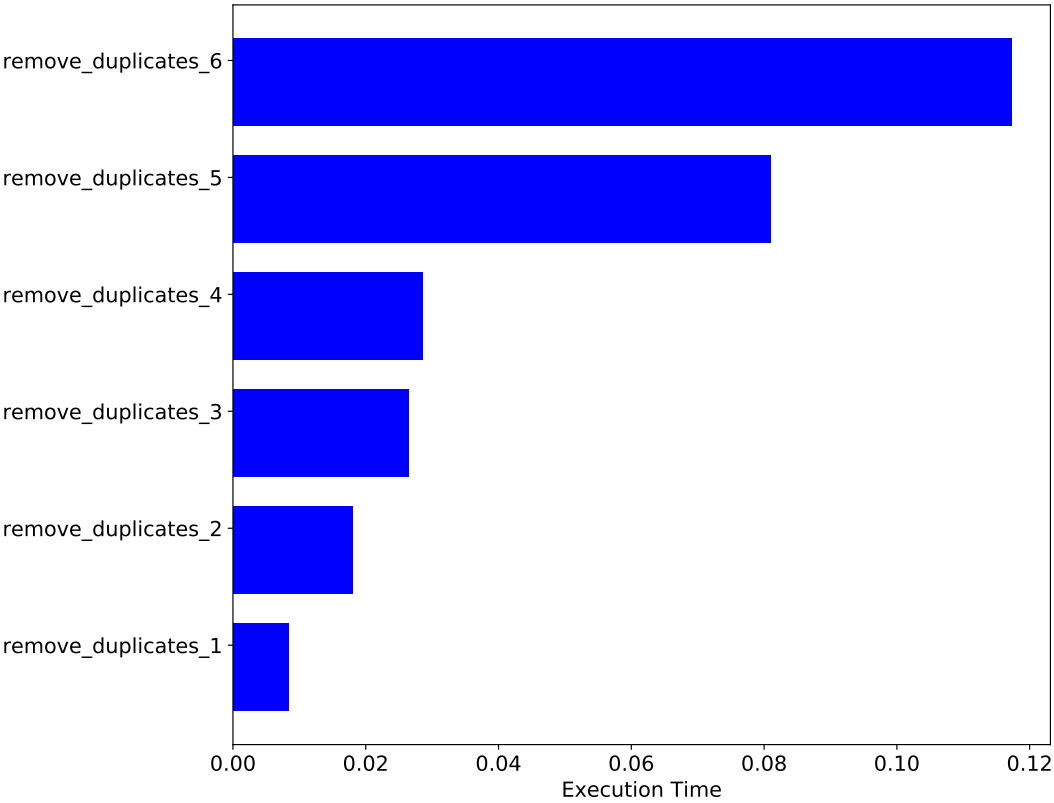
TIMES FOR THE FIRST 196000 NUMBERS



TIMES FOR THE FIRST 198000 NUMBERS



TIMES FOR THE FIRST 200000 NUMBERS



7. Bibliografía

Referencias

- [1] Michel Goossens, Frank Mittelbach, and Alexander Samarin. *The L^AT_EX Companion*. Addison-Wesley, Reading, Massachusetts, 1993.
- [2] Git. *Sobre el Control de Versiones*. [Una breve historia de Git]. <https://git-scm.com/book/es/v2/Inicio—Sobre-el-Control-de-Versiones-Una-breve-historia-de-Git>
- [3] Wikipedia *Análisis de rendimiento de software*. https://es.wikipedia.org/wiki/An%C3%A1lisis_de_rendimiento_de_software
- [4] Matplotlib *Generación de gráficos con Matplotlib* <https://matplotlib.org/>