

MyO – Trabajo Práctico

Ximena Ebertz

Hernán Rondelli

Lucía Soria

2^{do} Semestre 2024

Contenidos

1	Introducción	2
2	Parte 1	2
2.1	Descripción	2
2.2	Modelo Lineal	2
2.2.1	Definición de variables	3
2.2.2	Modelo	3
2.2.3	Restricción: que no se elijan discos vacíos	3
2.2.4	Modelo Zimpl	5
2.2.5	Modelo PySCIPOpt	6
3	Parte 2	13
3.1	Descripción	13
3.2	Modelo Lineal	13
3.2.1	Definición de variables	13
3.2.2	Modelo	14
3.2.3	Modelo Zimpl	15
3.2.4	Modelo PySCIPOpt	16
4	Parte 3	23
4.1	Descripción	23
4.2	Modelo Lineal	23
4.2.1	Definición de variables	23
4.2.2	Modelo	23
5	Parte 4	24
5.1	Descripción	24
5.2	Modelo Lineal	24
5.2.1	Definición de variables	24
5.2.2	Modelo	25

6	Parte 5	25
6.1	Descripción	25
6.2	Pseudocódigo de alto nivel	26
7	Parte 6	27
7.1	Descripción	27
8	Resultados de script	27
9	Conclusiones	29

1 Introducción

Este trabajo práctico tiene como objetivo aplicar conceptos de modelado y de optimización, utilizando técnicas de programación lineal, y distintas tecnologías y herramientas. Buscamos, a través de una serie de problemas relacionados con la gestión de archivos en una empresa de datos, diseñar soluciones eficientes que optimicen el uso de recursos, como el almacenamiento en discos y la importancia de los datos.

La actividad se divide en seis partes, cada una con enfoques específicos que abarcan desde el modelado directo en `SCIP`, hasta la implementación de algoritmos en Python con la biblioteca `PySCIP0pt`. A lo largo de este trabajo, se considerarán restricciones como la capacidad limitada de los discos, la importancia de los archivos, y la estructura de conjuntos para maximizar la eficiencia del almacenamiento y la relevancia de los datos seleccionados.

2 Parte 1

2.1 Descripción

Se deben realizar backups de n archivos de la empresa *BigData* en discos idénticos de capacidad d , sin dividir los archivos. Cada archivo es menor que la capacidad de un disco. Se necesita saber cuál es la cantidad mínima de discos necesarios para almacenar todos los archivos, considerando la capacidad de cada disco y el tamaño de cada archivo.

2.2 Modelo Lineal

Para minimizar la cantidad de discos, usamos una variable binaria y_j que nos indica si se está usando el disco j o no, y otra variable binaria x_{ij} para identificar si cada archivo i se guarda en el disco j . Luego, establecimos dos restricciones: la primera se asegura de que cada archivo esté en un solo disco, y la segunda, que los tamaños de los archivos de cada disco no superen la capacidad del disco.

2.2.1 Definición de variables

F : conjunto de nombres de cada archivo i , con $i = 1, \dots, n$.

s_i : tamaño del archivo i en MB, con $i = 1, \dots, n$.

n : cantidad de archivos.

m : cantidad de discos disponibles, a lo sumo, un archivo por disco, $m = n$.

D : conjunto de j discos de tamaño d , con $j = 1, \dots, m$.

d : tamaño del disco en TB.

x_{ij} binaria: 1 indica si el archivo i está en el disco j , $\forall i = \{1, \dots, n\}$, $\forall j = \{1, \dots, m\}$, donde $x_{ij} \in \{0, 1\}$, 0 en el caso de que no esté

y_j binaria: 1 indica que se está usando el disco j .

2.2.2 Modelo

$$\begin{aligned} & \text{Minimize} && \sum_{j=1}^m y_j \\ & \text{Subject to} && \sum_{j=1}^m x_{ij} = 1 \quad \forall i \in \{1, \dots, n\}, \text{ (el archivo } i \text{ debe estar en un solo disco)} \\ & && \sum_{i=1}^n s_i x_{ij} \leq d \cdot y_j \quad \forall j \in \{1, \dots, m\}, \text{ (los archivos que entran en el disco } j) \\ & && x_{ij} \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\}, \forall j \in \{1, \dots, m\} \\ & && y_j \in \{0, 1\}, \quad \forall j \in \{1, \dots, m\} \end{aligned}$$

2.2.3 Restricción: que no se elijan discos vacíos

Inicialmente, habíamos planteado el programa lineal de la Sección 2.2 con una tercera restricción:

$$\sum_{i=1}^n x_{ij} \leq n \cdot y_j \quad \forall j \in \{1, \dots, m\}, \text{ (no se puede usar un disco vacío)}$$

En esta sección, explicamos porqué esta restricción—aunque no modificaba el resultado de la función objetivo—no es necesaria, mediante una demostración por el absurdo.

Supongamos que contamos con dos archivos f_1 y f_2 , y dos discos, d_1 y d_2 . En este caso, una solución óptima hipotética sería $y = (1, 1)$, que indica que se utilizan los dos discos.

Supongamos que esta solución considera que los archivos f_1 y f_2 se almacenan en el disco d_2 (considerando que ambos archivos entran en el disco), quedando el disco d_1 vacío. Si analizamos las restricciones, vemos que todas se cumplen. Es decir:

$$\sum_{j=1}^2 x_{ij} = 1 \quad \forall i \in \{1, 2\} \quad (1)$$

$$\sum_{i=1}^2 s_i x_{ij} \leq d \cdot y_j \quad \forall j \in \{1, 2\} \quad (2)$$

Vemos que (1) se cumple ya que ambos archivos f_1 y f_2 están en un sólo disco, d_2 , por hipótesis. Por otro lado, (2) se cumple, ya que...

(i) $0 \leq d_1 \cdot 1$, ya que el disco está vacío.

(ii) $s_1 + s_2 \leq d_2 \cdot 1$, por hipótesis.

A parte, notar que $x_{ij} \in \{0, 1\}$, $\forall i \in \{1, 2\}$, $\forall j \in \{1, 2\}$, y $y_j \in \{0, 1\}$, $\forall j \in \{1, 2\}$ por hipótesis.

Queremos ver, a partir de esto, que existe una solución x^* mejor, por lo que x no sería una solución óptima.

Consideremos la solución $x^* = (0, 1)$. Vemos que también se cumplen las restricciones:

$$\sum_{j=1}^2 x_{ij} = 1 \quad \forall i \in \{1, 2\} \quad (3)$$

$$\sum_{i=1}^2 s_i x_{ij} \leq d \cdot y_j \quad \forall j \in \{1, 2\} \quad (4)$$

en (3) se cumple ya que ambos archivos f_1 y f_2 están en un sólo disco, d_2 , por hipótesis, y (4) se cumple, ya que...

(i) $0 \leq d_1 \cdot 0$, ya que no se elige el disco 1.

(ii) $s_1 + s_2 \leq d_2 \cdot 1$, por hipótesis.

También se cumplen las condiciones con respecto a los valores que pueden tomar $x_{1,j}$ e y_j , ya que no fueron modificadas. Entonces, x^* es una solución factible.

Para x ,

$$\sum_{j=1}^m y_j = 2.$$

Por otra parte, para x^* ,

$$\sum_{j=1}^m y_j = 1.$$

Entonces, x^* es una mejor solución que x , por lo que x no puede ser una solución óptima. Este absurdo surgió de considerar que puede haber un disco vacío en una solución óptima de nuestro problema lineal. Esto implica que la restricción planteada no aporta valor al problema, ya que probamos que se llega a una solución correcta sin ella.

2.2.4 Modelo Zimpl

A continuación presentamos el modelo Zimpl, con los comentarios correspondientes que indican la relación con el modelo lineal de la sección anterior.

```
1 param input := "IN/a_1.in"; # Path del archivo de input
2
3 # Tamaño del disco en terabytes
4 param d_t := read input as "1n" use 1 comment "#";
5 param d := d_t * 10**6;
6
7 # Conjunto de archivos f_{i} \forall i \in \{1, \dots, n\}
8 set F := { read input as "<1s>" skip 2 comment "#" };
9
10 param n := card(F); # Cantidad de archivos
11 param m := n; # Cantidad de discos, a lo sumo, un disco por archivo
12
13 # Conjunto de discos
14 set D := { 1 .. m };
15
16 # Tamaños de los archivos f_{i} \forall i \in \{1, \dots, n\}
17 param s[F] := read input as "<1s> 2n" skip 2 comment "#";
18
19 # 1 si el archivo $i$ está en el disco $j$, 0 en caso contrario
20 var x[<i, j> in F * D] binary;
21
22 # 1 si se usa el disco $j$, 0 en caso contrario
23 var y[D] binary;
24
25 minimize disks: sum<j> in D: y[j];
26
27 # Cada archivo debe estar en un solo disco
28 subto c1:
29     forall <i> in F:
30         sum<j> in D: x[i, j] == 1;
31
32 # Los archivos $i$ que entran en el disco $j$
33 subto c2:
34     forall <j> in D:
35         sum<i> in F: s[i] * x[i, j] <= d * y[j];
```

2.2.5 Modelo PySCIPOpt

Como parte del trabajo, incluimos una implementación del modelo utilizando Python, con la librería PySCIPOpt. La siguiente implementación corresponde a una función equivalente al código presentado con Zimpl.

```
1 import time
2 import sys
3 from pycscipopt import Model
4 from pycscipopt import quicksum
5 from pycscipopt import SCIP_PARAMSETTING
6
7 def distribuir_archivos_1(d_t: int, F: list[str], s: list[int], time_limit=420):
8     # model = crear_modelo_1(d_t, F, s, time_limit)
9     model, fake_y, fake_x = resolver_modelo_binario_1(d_t, F, s, time_limit)
10
11     try:
12         x, obj_x = obtener_solucion_primal_1(model)
13     except TypeError:
14         x, obj_x = None, None
15
16     sys.stderr.write(f"[Debugging] obj x {obj_x}\n")
17
18     solution = model.getBestSol()
19     status = model.getStatus()
20
21     if solution is not None and status in ["optimal", "feasible"]:
22         return [F, model, fake_y, fake_x, s]
23     else:
24         return None
25
26 def resolver_modelo_binario_1(d_t: int, F: list[str], s: list[int], time_limit=420):
27     model = Model("model_part_1")
28
29     d = d_t * 10**6
30
31     n = len(F)
32     m = n # no se puede tener más discos que archivos
33
34     #####
35     # BINARY
36     # x_{i, j} = 1 si se elige el archivo i para el disco j, 0 si no
37     #
38     x = {}
```

```

39 for i in range(n):
40     for j in range(m):
41         x[i, j] = model.addVar(name = f"x_{i}_{j}", vtype="BINARY")
42
43     # y_{jj} = 1 si se elige el disco j, 0 si no
44     y = [model.addVar(f"y_{j}", vtype="BINARY") for j in range(m)]
45     #
46     #####
47
48     # minimize disks:
49     model.setObjective(quicksum(y), sense="minimize")
50
51     # Que los archivos se elijan solo para un disco
52     for i in range(n):
53         model.addCons(quicksum(x[i, j] for j in range(m)) == 1)
54
55     # Que no se pasen de capacidad los discos
56     for j in range(m):
57         model.addCons(quicksum(x[i, j] * s[i] for i in range(n)) <= d * y[j])
58
59     # Configurar el límite de tiempo en el solver
60     model.setParam("limits/time", time_limit)
61     model.setParam("display/freq", 1)
62
63     model.optimize()
64
65     sys.stderr.write(f"[Debugging] [MODELO 1] Time: {model.getSolvingTime()}\n\n")
66     sys.stderr.write(f"[Debugging] [MODELO 1] Cantidad sols: {model.getNSols()}\n\n")
67
68     return model, y, x
69
70 # Crea el modelo relajado y lo devuelve optimizado
71 def crear_modelo_1(d_t: int, F: list[str], s: list[int], time_limit=420):
72     sys.stderr.write(f"[Debugging] [MODELO 1] Inicio\n\n")
73     model = Model("model_part_1")
74
75     d = d_t * 10**6
76
77     n = len(F)
78     m = n # no se puede tener más discos que archivos
79
80     #####
81     # CONTINUOUS

```

```

82 #  $x_{\{i, j\}} = 1$  si se elige el archivo  $i$  para el disco  $j$ , 0 si no
83 #
84 x = {}
85 for i in range(n):
86     for j in range(m):
87         x[i, j] = model.addVar(name = f"x_{i}_{j}", lb=0, ub=1, vtype="CONTINUOUS")
88
89 #  $y_{\{j\}} = 1$  si se elige el disco  $j$ , 0 si no
90 y = [model.addVar(f"y_{j}", lb=0, ub=1, vtype="CONTINUOUS") for j in range(m)]
91 #
92 #####
93
94 # minimize disks:
95 model.setObjective(quicksum(y), sense="minimize")
96
97 # Que los archivos se elijan solo para un disco
98 for i in range(n):
99     model.addCons(quicksum(x[i, j] for j in range(m)) == 1)
100
101 # Que no se pasen de capacidad los discos
102 for j in range(m):
103     model.addCons(quicksum(x[i, j] * s[i] for i in range(n)) <= d * y[j])
104
105 # Configurar el límite de tiempo en el solver
106 model.setParam("limits/time", time_limit)
107 model.setParam("display/freq", 1)
108
109 # Desactivación temporal de presolve
110 model.setPresolve(SCIP_PARAMSETTING.OFF)
111
112 #####
113 # Esto es la clave para que  $obj(dual) = obj(primal)$ 
114 # según documentación de pycscipy
115 #
116 model.setHeuristics(SCIP_PARAMSETTING.OFF)
117 model.disablePropagation()
118 #
119 #####
120
121 model.optimize()
122
123 sys.stderr.write(f"[Debugging] [MODELO 1] Time: {model.getSolvingTime()}\n\n")
124 sys.stderr.write(f"[Debugging] [MODELO 1] Cantidad sols: {model.getNSols()}\n\n")

```



```

125
126     return model
127
128 def obtener_solucion_primal_1(model):
129     sol = model.getBestSol()
130     status = model.getStatus()
131
132     if sol is not None and status in ["optimal", "feasible"]:
133         sys.stderr.write(f"[Debugging] {model.getStatus()}: {model.getBestSol()}\n\n")
134
135         # El nombre de variable x, acá es misleading, es x porque es el primal.
136         # Esta línea devuelve **todas** las variables del modelo, las  $x_i$ : archivo seleccionado y las
137         x = [v.getLPSol() for v in model.getVars(False)]
138
139         sys.stderr.write(f"[Debugging]: todas las variables del primal model_part_1 {x}\n\n")
140         return x, model.getObjVal()
141     else:
142         return None
143
144 def obtener_solucion_dual_1(model):
145     # No debería ser necesario si el modelo ya viene optimizado con el presolving off
146     # Asegurarse de apagar el presolving
147     # model.setPresolve(SCIP_PARAMSETTING.OFF)
148
149     # La longitud de y coincide con la cantidad de archivos del input
150     y = [model.getDualSolVal(c) for c in model.getConss(False)]
151
152     return y, quicksum(y)

```

Esta función recibe los datos leídos de un archivo de configuración creado con datos aleatorios. El siguiente código muestra las funciones necesarias para generar este archivo.

```
1 import os
2 import random
3
4 def generar_archivos(cant_archivos):
5     archivos = {}
6     contador = 1
7     for i in range(cant_archivos):
8         archivos["archivo" + str(contador)] = random.randint(1000000, 10000000)
9         contador += 1
10    return archivos
11
12
13 def generar_configuracion(nombre_archivo):
14     capacidad_discos = random.randint(1, 300)
15     cant_archivos = random.randint(400, 550)
16     archivos = generar_archivos(cant_archivos)
17
18     # ruta_in = os.path.join(os.path.dirname(__file__), "IN", nombre_archivo)
19     ruta_in = os.path.dirname(__file__) + "/IN/" + nombre_archivo
20     with open(ruta_in, "w") as f:
21
22         f.write(f"# Capacidad de discos en TB (= 1.000.000 MB)\n")
23         f.write(str(capacidad_discos) + "\n")
24
25         f.write(f"\n# Cantidad de archivos para backup\n")
26         f.write(str(cant_archivos) + "\n")
27
28         f.write(f"\n# Archivos: archivo_id, tamaño (MB) \n")
29         for archivo in archivos:
30             f.write(archivo + " " + str(archivos[archivo]) + "\n")
31
32
33 def leer_configuracion(nombre_archivo):
34     capacidad_disco = 0
35     nombres_archivos = []
36     tamaños_archivos = []
37
38     # ruta_in = os.path.join(os.path.dirname(__file__), "IN", nombre_archivo)
39     ruta_in = os.path.dirname(__file__) + "/IN/" + nombre_archivo
40     with open(ruta_in, "r") as f:
41         lineas = f.readlines()
```

```

42     capacidad_disco = int(lineas[1].strip())
43
44     for i in range(7, len(lineas)):
45         if lineas[i].strip():
46             archivo = lineas[i].split()
47             nombres_archivos.append(archivo[0])
48             tamaños_archivos.append(int(archivo[1]))
49     return capacidad_disco, nombres_archivos, tamaños_archivos

```

La solución incluye un script ejecutable que recibe como parámetro el nombre del archivo a generar, genera el archivo y calcula el resultado del modelo. Luego, genera el output correspondiente.

```

1  #!/usr/bin/env python3
2
3  import sys
4  from configuracion_1 import *
5  from generador_output_1 import *
6  from model_part_1 import *
7
8  if len(sys.argv) != 2:
9      print(f"Uso: {sys.argv[0]} nombre_archivo")
10     sys.exit(1)
11
12  archivo = sys.argv[1]
13  print(f"Utilizando {archivo}\n")
14
15  # generar_configuracion(archivo)
16
17  capacidad_disco, nombres_archivos, tamaños_archivos = leer_configuracion(f"{archivo}")
18
19  solucion = distribuir_archivos_1(capacidad_disco, nombres_archivos, tamaños_archivos, 420)
20
21  if solucion is not None:
22      generar_output(f"{archivo[:-3]}.out", solucion)
23  else:
24      generar_output_fallido(f"{archivo[:-3]}.out")

```

Utilizamos las siguientes funciones para generar el output.

```
1 import os
2
3 def generar_output(nombre_archivo, solucion): # solucion = [F, model, y, x, s]
4     F = solucion[0]
5     model = solucion[1]
6     y = solucion[2]
7     x = solucion[3]
8     s = solucion[4]
9
10    cant_archivos = len(F)
11    cant_discos = round(float(model.getObjVal()))
12
13    # La cantidad de discos disponibles es a lo sumo la cantidad de archivos
14    number_of_disks = cant_archivos
15
16    ruta_out = os.path.join(os.path.dirname(__file__), ".", "OUT",
17                             nombre_archivo)
18    with open(ruta_out, "w") as f:
19        f.write(
20            f"Para la configuración del archivo, {cant_discos} discos son suficientes.\n"
21        )
22        for j in range(number_of_disks):
23            if model.getVal(y[j]) == 0:
24                continue
25
26            archivos_en_disco = []
27            espacio_ocupado = 0
28
29            for i in range(cant_archivos):
30                if model.getVal(x[i, j]) == 0:
31                    continue
32                archivos_en_disco.append(f"{F[i]} {s[i]}")
33                espacio_ocupado = espacio_ocupado + s[i]
34
35            f.write(f"\nDisco {j+1}: {espacio_ocupado} MB\n")
36
37            for archivo in archivos_en_disco:
38                f.write(archivo + "\n")
39
40
41 def generar_output_fallido(nombre_archivo):
42     ruta_out = os.path.join(os.path.dirname(__file__), ".", "OUT",
```

```

43         nombre_archivo)
44     with open(ruta_out, "w") as f:
45         f.write(
46             f"No se ha encontrado solucion para la configuracion del archivo.\n"
47         )

```

3 Parte 2

3.1 Descripción

Se asigna un indicador de importancia I a cada archivo según características como sensibilidad, frecuencia de uso, y rareza. El objetivo es encontrar un subconjunto de archivos que quepan en un disco de tamaño d , y que la suma de indicadores de importancia sea máxima.

3.2 Modelo Lineal

Para maximizar la importancia de los archivos en un disco, usamos una variable binaria x_i que indica si el archivo i está o no en el disco, y la multiplicamos por la constante I_i que indica la importancia del archivo i . La restricción que utilizamos para este problema lineal, es que la suma de los tamaños de los archivos no superen el tamaño del disco.

3.2.1 Definición de variables

F : conjunto de nombres de i archivos, con $i = 1, \dots, n$.

s_i constante: tamaño del archivo i en MB, con $i = 1, \dots, n$.

I_i constante: indicador de importancia para el archivo i , con $i = 1, \dots, n$.

n constante: cantidad de archivos.

d constante: tamaño del disco en TB.

x_i binaria: 1 indica si el archivo i está en el disco, $\forall i = \{1, \dots, n\}$, donde $x_i \in \{0, 1\}$, 0 en el caso de que no esté

3.2.2 Modelo

$$\begin{array}{ll}\text{Maximize} & \sum_{i=1}^n I_i \cdot x_i \\ \text{Subject to} & \sum_{i=1}^n s_i x_i \leq d, \quad (\text{archivos que entran en el disco}) \\ & x_i \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\} \\ & s_i \geq 0, \quad \forall i \in \{1, \dots, n\} \\ & I_i \geq 0, \quad \forall i \in \{1, \dots, n\} \\ & n \geq 0 \\ & d \geq 0\end{array}$$

3.2.3 Modelo Zimpl

A continuación presentamos el modelo Zimpl, con los comentarios correspondientes que indican la relación con el modelo lineal de la sección anterior.

```
1 param input := "IN/b_1.in"; # Path del archivo de input
2
3 # Tamaño del disco en terabytes
4 param d_t := read input as "1n" use 1 comment "#";
5 param d := d_t * 10**6;
6
7 # Conjunto de archivos f_{i} \forall i \in \{1, \dots, n\}
8 set F := { read input as "<1s>" skip 2 comment "#" };
9
10 param n := card(F); # Cantidad de archivos
11
12 # Tamaños de los archivos f_{i} \forall i \in \{1, \dots, n\}
13 param s[F] := read input as "<1s> 2n" skip 2 comment "#";
14
15 # Importancias de los archivos
16 param importance[F] := read input as "<1s> 3n" skip 2 comment "#";
17
18 # 1 si el archivo $i$ está en el disco, 0 en caso contrario
19 var x[F] binary;
20
21 maximize importances: sum<i> in F: importance[i] * x[i];
22
23 # Los archivos $i$ que entran en el disco
24 subto c1:
25     sum<i> in F: s[i] * x[i] <= d;
```

3.2.4 Modelo PySCIPOpt

De la misma manera que en el modelo anterior, generamos el modelo, la configuración y los archivos de output utilizando Python.

El modelo es el siguiente.

```
1 import sys
2 from pycscipopt import Model
3 from pycscipopt import quicksum
4 from pycscipopt import SCIP_PARAMSETTING
5
6 # Model segunda parte
7 def distribuir_archivos_2(d_t: int, F: list[str], s: list[int], I: list[float]):
8     model, fake_x = resolver_modelo_binario_2(d_t, F, s, I, 420)
9     try:
10         x, obj_x = obtener_solucion_primal_2(model)
11     except TypeError:
12         x, obj_x = None, None
13
14     # sys.stderr.write(f"[Debugging] obj x {obj_x}\n")
15
16     solution = model.getBestSol()
17     status = model.getStatus()
18
19     if solution is not None and status in ["optimal", "feasible"]:
20         sys.stderr.write(f"[Debuggin] {status}: {solution}\n\n")
21         return [F, model, fake_x, I, s]
22     else:
23         return None
24
25 def resolver_modelo_binario_2(d_t: int, F: list[str], s: list[int], I: list[float], time_limit=420)
26     model = Model("model_part_2")
27
28     d = d_t * 10**6
29
30     n = len(F)
31
32     #####
33     # BINARY
34     # x_{i} = 1 si se elige el archivo i, 0 si no
35     x = [model.addVar(f"y_{i}", vtype="BINARY") for i in range(n)]
36     #####
37
38     # maximize importance:
```



```

39 model.setObjective(quicksum(x[i] * I[i] for i in range(n)), sense="maximize")
40
41 # los archivos elegidos deben entrar en el disco
42 model.addCons(quicksum(x[i] * s[i] for i in range(n)) <= d)
43
44 # Configurar el límite de tiempo en el solver
45 model.setParam("limits/time", time_limit)
46 model.setParam("display/freq", 1)
47
48 model.optimize()
49
50 sys.stderr.write(f"[Debugging] model_part_2 Time: {model.getSolvingTime()}\n\n")
51 sys.stderr.write(f"[Debugging] model_part_2 Cantidad sols: {model.getNSols()}\n\n")
52
53 return model, x
54
55 # Crea el modelo relajado y lo devuelve optimizado
56 def crear_modelo_2(d_t: int, F: list[str], s: list[int], I: list[float], time_limit=420):
57     model = Model("model_part_2")
58
59     d = d_t * 10**6
60
61     n = len(F)
62
63     #####
64     # CONTINUOUS
65     # x_{i} = 1 si se elige el archivo i, 0 si no
66     x = [model.addVar(f"y_{i}", lb=0, ub=1, vtype="CONTINUOUS") for i in range(n)]
67     #####
68
69     # maximize importance:
70     model.setObjective(quicksum(x[i] * I[i] for i in range(n)), sense="maximize")
71
72     # los archivos elegidos deben entrar en el disco
73     model.addCons(quicksum(x[i] * s[i] for i in range(n)) <= d)
74
75     # Configurar el límite de tiempo en el solver
76     model.setParam("limits/time", time_limit)
77     model.setParam("display/freq", 1)
78
79     # Desactivación temporal de presolve
80     model.setPresolve(SCIP_PARAMSETTING.OFF)
81

```

```

82 #####
83 # Esto es la clave para que obj(dual) = obj(primal)
84 # según documentación de pycipopt
85 #
86 model.disablePropagation()
87 model.setHeuristics(SCIP_PARAMSETTING.OFF)
88 #
89 #####
90
91 model.optimize()
92
93 sys.stderr.write(f"[Debugging] model_part_2 Time: {model.getSolvingTime()}\n\n")
94 sys.stderr.write(f"[Debugging] model_part_2 Cantidad sols: {model.getNSols()}\n\n")
95
96 return model
97
98 def obtener_solucion_primal_2(model):
99     sol = model.getBestSol()
100     status = model.getStatus()
101
102     if sol is not None and status in ["optimal", "feasible"]:
103         sys.stderr.write(f"[Debugging] {status}: {sol}\n\n")
104
105         x = [v.getLPSol() for v in model.getVars(False)]
106
107         sys.stderr.write(f"[Debugging]: {x}\n\n")
108         return x, model.getObjVal()
109     else:
110         return None
111
112 def obtener_solucion_dual_2(model):
113     # No debería ser necesario si el modelo ya viene optimizado con el presolving off
114     # Asegurarse de apagar el presolving
115     model.setPresolve(SCIP_PARAMSETTING.OFF)
116
117     # La longitud de y coincide con la cantidad de archivos del input
118     y = [model.getDualSolVal(c) for c in model.getConss(False)]
119
120     return y, quicksum(y)

```

Los archivos de configuración se generan de manera similar al modelo previo.

```
1 import os
2 import random
3
4
5 def generar_archivos(cant_archivos):
6     archivos = {}
7     contador = 1
8     for i in range(cant_archivos):
9         # {nombreArchivo: [tamaño, importancia] }
10        archivos["archivo" + str(contador)] = [
11            random.randint(1000000, 10000000),
12            random.randint(1, 10)
13        ]
14        contador += 1
15    return archivos
16
17
18 def generar_configuracion(nombre_archivo):
19     capacidad_discos = random.randint(1, 100)
20     cant_archivos = random.randint(1, 50)
21     archivos = generar_archivos(cant_archivos)
22
23     ruta_in = os.path.join(os.path.dirname(__file__), ".", "IN",
24                             nombre_archivo)
25     with open(ruta_in, "w") as f:
26
27         f.write(f"# Capacidad de discos en TB (= 1.000.000 MB)\n")
28         f.write(str(capacidad_discos) + "\n")
29
30         f.write(f"\n# Cantidad de archivos para backup\n")
31         f.write(str(cant_archivos) + "\n")
32
33         f.write(f"\n# Archivos: archivo_id, tamaño (MB), importancia\n")
34         for archivo in archivos:
35             f.write(archivo + " " + str(archivos[archivo][0]) + " " +
36                     str(archivos[archivo][1]) + "\n")
37
38
39 def leer_configuracion(nombre_archivo):
40     capacidad_disco = 0
41     nombres_archivos = []
42     tamaños_archivos = []
```

```

43 importancias_archivos = []
44
45 ruta_in = os.path.join(os.path.dirname(__file__), ".", "IN",
46                         nombre_archivo)
47 with open(ruta_in, "r") as f:
48
49     lineas = f.readlines()
50     capacidad_disco = int(lineas[1].strip())
51     for i in range(7, len(lineas)):
52         if lineas[i].strip():
53             archivo = lineas[i].split()
54             nombres_archivos.append(archivo[0])
55             tamaños_archivos.append(int(archivo[1]))
56             importancias_archivos.append(int(archivo[2]))
57 return capacidad_disco, nombres_archivos, tamaños_archivos, importancias_archivos

```

El script ejecutable realiza los mismos pasos que el del modelo anterior, pero, considerando la importancia de cada archivo.

```
1 #!/usr/bin/env python3
2
3 import sys
4 from configuracion_2 import *
5 from generador_output_2 import *
6 from model_part_2 import *
7
8 if len(sys.argv) != 2:
9     print(f"Uso: {sys.argv[0]} nombre_archivo")
10    sys.exit(1)
11
12 archivo = sys.argv[1]
13 print(f"Utilizando {archivo}\n")
14
15 generar_configuracion(archivo)
16 capacidad_disco, nombres_archivos, tamaños_archivos, importancia_archivos = leer_configuracion(f"{a
17 solucion = distribuir_archivos_2(capacidad_disco, nombres_archivos, tamaños_archivos, importancia_a
18
19 if solucion is not None:
20     generar_output(f"{archivo[:-3]}.out", solucion)
21 else:
22     generar_output_fallido(f"{archivo[:-3]}.out")
```

Luego, el output se almacena de la siguiente manera.

```
1 import os
2
3 # solucion = F, model, x, I, s
4 def generar_output(nombre_archivo, solucion):
5     F = solucion[0]
6     model = solucion[1]
7     x = solucion[2]
8     I = solucion[3]
9     s = solucion[4]
10
11     cant_archivos = len(F)
12     cant_archivos_elegidos = sum(1 for i in range(cant_archivos)
13                                 if model.getVal(x[i]) > 0.5)
14     archivos_elegidos = []
15     importancia_total = 0
16
17     ruta_out = os.path.join(os.path.dirname(__file__), ".", "OUT",
18                             nombre_archivo)
19     with open(ruta_out, "w") as f:
20         f.write(
21             f"Para la configuracion del archivo, se han elegido {cant_archivos_elegidos} archivos.\n"
22         )
23
24         for i in range(cant_archivos):
25             if model.getVal(x[i]) > 0.5: # se eligio el archivo
26                 archivos_elegidos.append(f"{F[i]} {s[i]} {I[i]}")
27                 importancia_total += I[i]
28
29         for archivo in archivos_elegidos:
30             f.write(archivo + "\n")
31
32         f.write(
33             f"\nLa suma de sus indicadores de importancia da {importancia_total}."
34         )
35
36 def generar_output_fallido(nombre_archivo):
37     with open(nombre_archivo, "w") as f:
38         f.write(
39             f"No se ha encontrado solucion para la configuracion del archivo.\n"
40         )
```

4 Parte 3

4.1 Descripción

Dado un conjunto de archivos C y una familia de conjuntos de archivos \mathcal{H} , donde cada conjunto de archivos H de \mathcal{H} está incluido en C , y cada archivo de C pertenece a algún conjunto de la familia \mathcal{H} , se busca un subconjunto mínimo \mathcal{I} de \mathcal{H} tal que todos los archivos de C estén cubiertos por al menos un conjunto de \mathcal{I} .

4.2 Modelo Lineal

4.2.1 Definición de variables

n : cantidad de archivos.

m : cantidad de conjuntos H .

x_j binaria: 1 indica que el conjunto H_j pertenece a \mathcal{I} .

y_{ij} **constante** binaria: 1 indica que el conjunto H_j contiene el archivo i .

4.2.2 Modelo

$$\begin{aligned} & \text{Minimize} && \sum_{j=1}^m x_j \\ & \text{Subject to} && \sum_{j=1}^m y_{i,j} \cdot x_j \geq 1 \quad \forall i, \text{ (todos los archivos deben estar en al menos un conjunto elegido)} \\ & && x_j \in \{0, 1\}, \quad \forall j \in \{1, \dots, m\} \\ & && y_{i,j} \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\}, \quad \forall j \in \{1, \dots, m\} \\ & && n \geq 0, \\ & && m \geq 0, \end{aligned}$$

5 Parte 4

5.1 Descripción

Se deben realizar backups de n archivos en discos idénticos de capacidad d . Se sabe que los archivos son de pocos tamaños distintos, y que cada archivo es menor que la capacidad de un disco. Se necesita saber cuál es la cantidad mínima de discos necesarios para almacenar todos los archivos, considerando la capacidad de cada disco y los tamaños de archivos.

5.2 Modelo Lineal

5.2.1 Definición de variables

F : conjunto de nombres de cada archivo i , con $i \in \{1, \dots, n\}$.

s_i : tamaño del archivo i en MB, con $i \in \{1, \dots, n\}$.

n : cantidad de archivos.

m : cantidad de discos disponibles, a lo sumo, un archivo por disco, $m = n$.

D : conjunto de j discos de tamaño d , con $j \in \{1, \dots, m\}$.

S : conjunto de los tamaños de s_i , con $i \in \{1, \dots, n\}$.

d : tamaño del disco en TB.

t : constante, máxima cantidad de tamaños que entran en un disco.

x_{ij} : binaria, toma el valor 1 si el archivo i está en el disco j , y 0 en caso contrario, para $\forall i \in \{1, \dots, n\}$ y $\forall j \in \{1, \dots, m\}$, donde $x_{ij} \in \{0, 1\}$.

y_j : binaria, toma el valor 1 si se está usando el disco j , y 0 en caso contrario, para $\forall j \in \{1, \dots, m\}$, donde $y_j \in \{0, 1\}$.

z_{sj} : binaria, toma el valor 1 si el tamaño s está en el disco j al menos una vez, y 0 en caso contrario, para $\forall s$ y $\forall j \in \{1, \dots, m\}$, donde $z_{sj} \in \{0, 1\}$.

M_{si} : constante, toma la cantidad de archivos de tamaño s_i que entran en un disco de tamaño d . Nótese que esto es equivalente a $\text{int}(toMB(d)/s)$.

5.2.2 Modelo

$$\begin{aligned}
& \text{Minimize} && \sum_{j=1}^m y_j \\
& \text{Subject to} && \sum_{j=1}^m x_{ij} = 1 \quad \forall i \in \{1, \dots, n\}, \text{ (el archivo } i \text{ debe estar en un solo disco)} \\
& && \sum_{i=1}^n s_i x_{ij} \leq d \cdot y_j \quad \forall j \in \{1, \dots, m\} \text{ (los archivos que entran en el disco } j \text{)} \\
& && z_{sj} \cdot M_{si} \leq \sum_{i=1}^n x_{ij} \quad \text{si } S_i = s, \quad \forall s \in S, \forall j \in \{1, \dots, m\}, \forall i \in \{1, \dots, n\} \text{ (marcamos los tamaños } c \text{)} \\
& && \sum_{s \in S} z_{sj} \leq t, \quad \forall j \in \{1, \dots, m\} \text{ (en cada disco no entran mas de } t \text{ tamaños)} \\
& && x_{ij} \in \{0, 1\}, \quad \forall i \in \{1, \dots, n\}, \forall j \in \{1, \dots, m\} \\
& && y_j \in \{0, 1\}, \quad \forall j \in \{1, \dots, m\} \\
& && z_{sj} \in \{0, 1\}, \quad \forall s \in S, \forall j \in \{1, \dots, m\} \\
& && n \geq 0, n \in \mathbb{Z} \\
& && s_i \geq 0, \forall i \in \{1, \dots, n\} \\
& && m \geq 0, m \leq n, m \in \mathbb{Z} \\
& && d \geq 0 \\
& && t \geq 0, t \in \mathbb{Z}
\end{aligned}$$

6 Parte 5

6.1 Descripción

El objetivo de este ejercicio es resolver el problema de las Secciones 2 y 5 suponiendo que los tamaños de los archivos pertenecen a un conjunto pequeño. Para ello, se genera una familia de conjuntos H que cubra todos los archivos y quepa en los discos, se plantea un modelo de programación lineal relajado P y se resuelve su dual para asignar valores a los archivos. Usando estos valores como indicadores de importancia, se selecciona un subconjunto de archivos. Si la solución es mejorable, se actualiza H y se repite el proceso, ajustando finalmente la solución para que sea entera.

6.2 Pseudocódigo de alto nivel

```
1 INICIO obtener_conjuntos(archivo, threshold = infinito)
2   capacidad_disco, nombres_archivos, tamanos_archivos <- leer_configuracion(archivo)
3
4   // PASO 1: Generar los conjuntos iniciales
5   conjuntos <- generar_conjuntos(capacidad_disco * 10^6, nombres_archivos, tamanos_archivos)
6
7   duracion <- threshold * 60
8   tiempo_inicio <- tiempo_actual()
9
10  MIENTRAS VERDADERO
11    // PASO 2 y 3: Resolver el modelo y obtener soluciones primal y dual
12    modelo <- crear_modelo_3(nombres_archivos, conjuntos)
13
14    x, obj_x <- obtener_solucion_primal_3(modelo)
15    y, obj_y <- obtener_solucion_dual_3(modelo)
16
17    optimo <- es_optimo(modelo, x)
18
19    // PASO 4: Distribuir archivos y generar solucion
20    distribucion <- distribuir_archivos_2(capacidad_disco, nombres_archivos, tamanos_archivos, y)
21    solucion_modelo_2 <- generar_output_modelo_2(distribucion)
22
23    // PASO 5: Verificar y actualizar conjuntos
24    SI suma(solucion_modelo_2[1]) > 1 Y tiempo_actual() - tiempo_inicio <= duracion ENTONCES
25      AGREGAR conjunto(solucion_modelo_2[0]) a conjuntos
26    SINO
27      IMPRIMIR suma(solucion_modelo_2[1])
28      SALIR DEL BUCLE
29    FIN SI
30  FIN MIENTRAS
31
32  x_estrella_int <- obtener_solucion_entera(modelo, x)
33  RETORNAR modelo
34 FIN
```

7 Parte 6

7.1 Descripción

En esta sección, se explican mejoras al algoritmo de la Sección 6.

Para agilizar el algoritmo, investigamos en la documentación de PySCIPOpt con el objetivo de identificar parámetros y configuraciones que permitan optimizar el tiempo de ejecución del código. PySCIPOpt, al estar basado en el solucionador SCIP, brinda muchas opciones para personalizar su comportamiento, como ajustes en las estrategias de búsqueda, límites en la exploración del árbol y técnicas específicas para encontrar soluciones factibles de manera más eficiente. Estas configuraciones permiten adaptar el rendimiento del solucionador a las características particulares del problema, maximizando su eficacia y reduciendo el tiempo computacional necesario. En nuestro caso, realizamos las siguientes mejoras al modelo:

8 Resultados de script

caso	cant	cota dual	mejor1	var1	tiempo1	mejor4	var4	tiempo4	mejor5	var5	tiempo5	mejor6
f0015.in	15	8.0	8.0	240	0.0	8.0	240	0.0	7.666666666666667	43	0.0	0
f0018.in	18	10.0	8.0	240	0.0	9.999999544151857	342	98.0	9.166666666666666	69	0.0	0
f0032.in	32	16.5	8.0	240	0.0	9.999999544151857	342	98.0	16.5	285	0.0	0

9 Conclusiones

Logramos familiarizarnos con el uso de SCIP, y con PySCIPOpt. Pudimos apreciar el alto nivel de abstracción de la sintaxis de Zimpl, utilizado directamente SCIP.

Este trabajo nos ayudó a familiarizarnos con las herramientas de modelado, y pudimos relacionarlas y aplicar lo visto en la materia.

Referencias

- [1] Koch, T. (2024) *Zimpl User Guide*, Berlin: ZIB.