

1. Introducción

En esta práctica he implementado en python el algoritmo de codificación de Huffman, uno de decodificación y funciones auxiliares para comprobar el Primer Teorema de Shannon.

2. Método

He implementado los algoritmos mediante programación orientada a objetos para aprovechar ciertas comodidades de python como la sobrecarga de operadores. Gracias a esto tambien se obtiene un código bastante modular y reutilizable.

Además he implementado un método `graph` para visualizar los árboles de huffman (he incluido ejemplos de visualizaciones en los anexos).

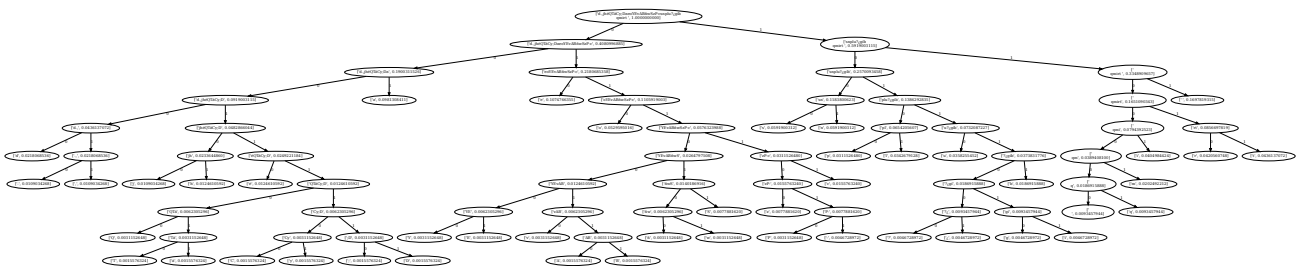
Para comprobar el Primer Teorema de Shannon he utilizado la siguiente definición de entropía, para un sistema de estados no equiprobables:

$$H(X) = - \sum_{i=1}^n P(x_i) \log_2 P(x_i)$$

donde X es una variable aleatoria discreta con x_1, \dots, x_n posibles estados que ocurren con probabilidades $P(x_1), \dots, P(x_n)$. En nuestro caso $C = X$ es el conjunto de caracteres de un idioma y las probabilidades son las obtenidas a partir de las frecuencias de las muestras (los textos anexos).

3. Resultados

En la siguiente gráfica (incluida como vectorial para poder visualizarse ampliando el documento) se puede ver el árbol de Huffman correspondiente al texto de referencia en Castellano.



3.1. Comprobación del Primer Teorema de Shannon

El Primer Teorema de Shannon nos da la siguiente cota para la longitud esperada de la codificación de Huffman:

$$H(C) \leq L(C) < H(C) + 1$$

donde $H(C)$ es la entropía del sistema y $L(C)$ la longitud esperada.

Los valores obtenidos para las muestras de textos en Español e Inglés son:

- Español: $H(C) = 4.4074$ y $L(C) = 4.4439$ con lo que

$$H(C) = 4.4074 \leq L(C) = 4.4439 < H(C) + 1 = 5.4074.$$

- Inglés: $H(C) = 4.1309$ y $L(C) = 4.1726$ con lo que

$$H(C) = 4.1309 \leq L(C) = 4.1726 < H(C) + 1 = 5.1309.$$

3.2. Codificación de la palabra 'medieval'

Una posible codificación binaria no comprimida para la palabra 'medieval' es la dada por la codificación *ASCII*, que codifica cada caracter mediante un byte (ocho bits). Por tanto la longitud de la codificación 'medieval' es $8 \cdot \text{long}(\text{'medieval'}) = 8 \cdot 8 = 64$.

Utilizando una codificación de Huffman se consigue utilizar códigos más cortos para los caracteres de mayor frecuencia, con la finalidad de reducir la longitud final de la codificación. Estas son las codificaciones que he obtenido para la palabra medieval, con sus longitudes correspondientes:

- Español: El código obtenido es 11000101000000110010100111001000110101, de longitud 38.
- Inglés: El código obtenido es 010011101110111110011000111010100000111001110101, de longitud 48.

Con lo que se comprueba que en cierto sentido se ha conseguido el objetivo de reducir las codificaciones. El punto delicado es que las frecuencias esperadas de caracteres con las que se construyen los árboles de Huffman se correspondan realmente con las frecuencias de los textos que queremos codificar. De no ser así podríamos generar codificaciones incluso menos óptimas que la *ASCII*.

3.3. Decodificación de un código dado

He hecho una implementación del algoritmo de Huffman inspirada en la implementación clásica, utilizando montículos. Debido a esto el árbol resultante es distinto al de la plantilla, puesto que el algoritmo de ordenación de montículos no tiene la misma estabilidad que el sort de numpy.

Por tanto con mi implementación no es posible decodificar el código proporcionado en el enunciado de la entrega. En lugar de esto he incluido en el código comprobaciones de que los códigos de la palabra 'medieval' obtenidos en el apartado anterior se decodifican correctamente en la palabra 'medieval'.

4. Código

El siguiente código con la implementación también está adjunto en la entrega y disponible, junto con esta memoria, en un repositorio git en el siguiente enlace: github.com/haztecaso/gcomp22.

```
4  """
5  Autor: Adrián Lattes Grassi
6  """
7
8  from decimal import Decimal
9  from functools import reduce
10 from heapq import heappop, heappush
11 from os import remove as remove_file
12 from typing import Dict, List, Optional
13 from math import log2
14
15
16 def frecuencias(texto:str)-> Dict[str, Decimal]:
17     """
18     Dado un texto calcula las frecuencias de cada caracter del texto.
19     """
20     letras = dict()
21     for char in texto:
22         if char not in letras:
23             letras[char] = Decimal(1)
24         else:
25             letras[char] += 1
26     n = len(texto)
27     return {char:frec/n for (char, frec) in letras.items()}
28
29
30 classCodigo():
31     """
32     Clase para códigos binarios.
33     """
34     def __init__(self, codigo:List[bool] = None):
35         self._codigo = codigo if codigo is not None else []
36
37     def pre(self, valor:bool):
38         """
39         Añade un valor al inicio del código
40         """
41         return Codigo([valor]+self._codigo)
42
```

```

43 @property
44 def vacio(self):
45     return len(self._codigo) == 0
46
47 def __add__(self, other):
48     """
49     Sobrecarga del operador + para concatenar códigos
50     """
51     returnCodigo(self._codigo + other._codigo)
52
53 def __len__(self):
54     return len(self._codigo)
55
56 def __iter__(self):
57     yield from self._codigo
58
59 def __repr__(self):
60     return f"Código<{''.join(map(lambda b: '1' if b else '0', self._codigo))>"
61
62
63 class ArbolHuffman():
64     """
65     Clase para árboles de Huffman.
66     """
67     def __init__(self, **kwargs):
68         assert 'clave' in kwargs or ('iz' in kwargs and 'dr' in kwargs)
69         if 'clave' in kwargs:
70             self.hoja = True
71             assert 'peso' in kwargs
72             self.peso = kwargs['peso']
73             self.clave = kwargs['clave']
74         else:
75             self.hoja = False
76             self.iz = kwargs['iz']
77             self.dr = kwargs['dr']
78             self.peso = self.iz.peso + self.dr.peso
79             self.clave = self.iz.clave + self.dr.clave
80             self._tabla_codigos = None
81
82     def __lt__(self, other):
83         return self.peso < other.peso
84
85     def __eq__(self, other):
86         return self.peso == other.peso
87
88     def __repr__(self):
89         return f"[{'self.clave'}, {self.peso:.10f}]"
90
91 @property
92 def tabla_codigos(self):
93     if self._tabla_codigos is None:
94         if self.hoja:
95             self._tabla_codigos = {self.clave:Codigo()}
96         else:
97             codigos_iz = {clave:codigo.pre(False) for (clave, codigo) in self.iz.tabla_codigos.items()}
98             codigos_dr = {clave:codigo.pre(True) for (clave, codigo) in self.dr.tabla_codigos.items()}
99             self._tabla_codigos:Optional[Dict[str,Codigo]] = {**codigos_iz, **codigos_dr}
100     return self._tabla_codigos
101
102 def codificar(self, data:str):
103     """
104     Codifica una cadena de caracteres
105     """
106     return reduce(lambda x,y: x+y,
107                  map(lambda e: self.tabla_codigos[e], data),
108                  Codigo())
109
110 def decodificar(self, codigo:Codigo):
111     """
112     Decodifica un código (o secuencia de códigos) de Huffman.
113     """
114     actual = self
115     codigo_actual = Codigo()
116     result = ""
117     for b in codigo:
118         actual = actual.dr if b else actual.iz
119         codigo_actual += Codigo([b])
120         if actual.hoja:
121             result += actual.clave
122             actual = self
123             codigo_actual = Codigo()
124     assert codigo_actual.vacio, f"No se ha terminado de decodificar el código. Código restante: {codigo_actual}"
125     return result
126
127 def graph(self, dot = None, render:bool = True, title:str='arbol'):
128     """
129     Exporta un dibujo del árbol en pdf
130     """
131     try:
132         from graphviz import Digraph
133     except ModuleNotFoundError as e:

```

```

134         print(f"ATENCIÓN: Para generar la gráfica de un ArbolHuffman es necesario instalar el paquete {e.name}.")
135     else:
136         if dot is None:
137             dot = Digraph(comment = title)
138         if self.hoja:
139             dot.node(self.clave, repr(self))
140         else:
141             dot.node(self.clave, repr(self))
142             self.iz.graph(dot, False)
143             self.dr.graph(dot, False)
144             dot.edge(self.clave, self.iz.clave, label="0")
145             dot.edge(self.clave, self.dr.clave, label="1")
146         if render:
147             print(f"Árbol de Huffman guardado en {title}.pdf")
148             dot.render(title)
149             remove_file(title) # Borrando archivo extra que crea graphviz
150         return dot
151
152
153 def huffman(frecs:Dict[str, Decimal]) -> ArbolHuffman:
154     """
155     Implementación del algoritmo de Huffman.
156     """
157     heap = []
158     for clave, peso in frecs.items():
159         heappush(heap, ArbolHuffman(clave = clave, peso = peso))
160     while len(heap) > 1:
161         iz = heappop(heap)
162         dr = heappop(heap)
163         heappush(heap, ArbolHuffman(iz = iz, dr = dr))
164     return heap[0]
165
166
167 def longitud_media(frecuencias:Dict[str, Decimal], tabla_codigos:Dict[str,Codigo]) -> Decimal:
168     """
169     Calcula la longitud esperada de la codificación de Huffman de un caracter,
170     dadas las frecuencias de los caracteres y la tabla de códigos de Huffman.
171     """
172     result = Decimal(0)
173     for clave, peso in frecuencias.items():
174         result += peso*len(tabla_codigos[clave])
175     return result
176
177
178 def entropia(frecuencias:Dict[str, Decimal]):
179     """
180     Entropía de Shannon calculada según la definición de wikipedia:
181     https://en.wikipedia.org/wiki/Entropy_(information_theory)
182     """
183     result = Decimal(0)
184     for peso in frecuencias.values():
185         result += peso*Decimal(log2(peso))
186     return -result
187
188
189 def main():
190     with open("GCOM2022_pract2_auxiliar_esp.txt") as f:
191         texto = '\n'.join(f.readlines())
192         frec_es = frecuencias(texto)
193     with open("GCOM2022_pract2_auxiliar_eng.txt") as f:
194         texto = '\n'.join(f.readlines())
195         frec_en = frecuencias(texto)
196
197     arbol_es = huffman(frec_es)
198     arbol_en = huffman(frec_en)
199     arbol_es.graph(title = 'arbol_es')
200     arbol_en.graph(title = 'arbol_en')
201     print()
202
203     print(f"""i) Comprobación del primer teorema de Shannon:
204     - Español: H(C) = {entropia(frec_es):.4f}, L(C) = {longitud_media(frec_es, arbol_es.tabla_codigos):.4f}
205     - Inglés: H(C) = {entropia(frec_en):.4f}, L(C) = {longitud_media(frec_en, arbol_en.tabla_codigos):.4f}
206     """)
207
208     medieval = "medieval"
209     codigo_medieval_es = arbol_es.codificar(medieval)
210     codigo_medieval_en = arbol_en.codificar(medieval)
211     print(f"""ii) Codificación de la palabra '{medieval}' en los dos idiomas:
212     - Español: {codigo_medieval_es} (longitud {len(codigo_medieval_es)})
213     - Inglés: {codigo_medieval_en} (longitud {len(codigo_medieval_en)})
214     - Con la codificación ASCII se utiliza un número fijo de bytes igual a 8,
215       con lo que la longitud de la palabra {medieval} sería {len(medieval)*8}.
216     """)
217
218     codigo = Codigo(list(map(lambda d:True if d == '1' else False, "101111011011101110111111")))
219     print(f"iii A) Decodificación del código {codigo}:")
220     try:
221         print(f"{codigo, arbol_en.decodificar(codigo) = }\n")
222     except AssertionError:
223         print(f"ERROR: No se ha podido decodificar el código {codigo}\n")

```

```
224 |
225 |     print("iii B) Comprobación de que la decodificación funciona correctamente:")
226 |     print(f"    - Español: {arbol_es.decodificar(codigo_medieval_es) = }")
227 |     print(f"    - Inglés: {arbol_en.decodificar(codigo_medieval_en) = }")
228 |
229 |
230 | if __name__ == "__main__":
231 |     main()
```