

```

# -*- coding: windows-1251 -*-
import string
import math
import heapq
from collections import Counter, namedtuple

f = open(r"res/1.txt", "r", encoding="windows-1251")
text = f.read()
text = text.lower()

def remove_chars_from_text(txt, chars):
    return ''.join([ch for ch in txt if ch not in chars])

spec_chars = string.punctuation + '«»-...1234567890—
qwertyuiopasdfghjklzxcvbnm\n\t\xa0©'
text = remove_chars_from_text(text, spec_chars)
text = text.replace('ë', 'e')
text = text.replace('Ъ', 'Ь')

text_tokens = dict()

def shannon(txt, n):
    for d in range(1, n):
        for i in range(len(txt) - d + 1):
            if text_tokens.get(txt[i:i + d]) is None:
                text_tokens[txt[i:i + d]] = 1
            else:
                text_tokens[txt[i:i + d]] = text_tokens.get(txt[i:i + d]) + 1
    if d == 1:
        max_h = math.log(len(text_tokens), 2)
        print("max_h = ", max_h)
        h = 0
        summa = sum(text_tokens.values())
        for value in text_tokens.values():
            h += (-1) * (value / summa) * math.log(value / summa, 2)
        h /= d
        print(f"h{d} = {h}")
        text_tokens.clear()
    # return h

class Node(namedtuple("Node", ["left", "right"])):
    def walk(self, code, acc):
        self.left.walk(code, acc + "0")
        self.right.walk(code, acc + "1")

```

```

class Leaf(namedtuple("Leaf", ["char"])):
    def walk(self, code, acc):
        code[self.char] = acc or "0"

def huffman_encode(s):
    h = []
    for ch, freq in Counter(s).items():
        h.append((freq, len(h), Leaf(ch)))
    heapq.heapify(h)

    count = len(h)
    while len(h) > 1:
        freq1, _count1, left = heapq.heappop(h)
        freq2, _count2, right = heapq.heappop(h)
        heapq.heappush(h, (freq1 + freq2, count, Node(left, right)))
        count += 1
    code = {}
    if h:
        [(_freq, _count, root)] = h
        root.walk(code, "")
    return code

def shannon1(txt):
    tokens = Counter(txt)
    h = 0
    summa = sum(tokens.values())
    for value in tokens.values():
        h += (-1) * (value / summa) * math.log(value / summa, 2)
    return h

def main():
    code = huffman_encode(text)
    # for ch in sorted(code.items(), key=lambda item: len(item[1])):
    #     print("{}: {}".format(ch[0], ch[1]))
    encode_text = "".join(code[ch] for ch in text)
    tokens = Counter(text)
    summa = sum(tokens.values())
    sr_dl = 0
    for ch in sorted(code):
        sr_dl += tokens[ch]/summa * len(code[ch])
    shannon(encode_text, 4)
    print("Средняя длина {}".format(sr_dl))
    print("Избыточность кодирования {}".format(sr_dl - 2))

```

main()