

```

1 #!python
2 # -*- coding:utf8 -*-
3
4 import re, sys
5
6 num = 0
7 def spannumber(m):
8     global num
9     span = m.group(1)
10    if span == 'span':
11        num += 1
12        return '<span' + str(num)
13    else:
14        temp = '</span' + str(num)
15        num -= 1
16        return temp
17
18 file_name = sys.argv[1]
19 input_file = open(file_name, "rU")
20 all = input_file.read()
21 input_file.close()
22
23 all = unicode(all, "sjis")
24
25 re_kaigyo = re.compile(r'¥n{2,}')
26 all = re_kaigyo.sub('{¥kaigyou}', all)
27 all = all.replace("¥n", " ")
28 all = all.replace("¥kaigyou", "¥n¥n")
29
30 re_generator = re.compile(r'<(meta name|link).+?>')
31 all = re_generator.sub('', all)
32
33 all = all.replace('<style>', '<meta http-equiv="content-style-type" content="text/css">')
34
35 all = all.replace('</style>', '<link rel="stylesheet" type="text/css" href="style.css">')
36
37 re_two_spaces = re.compile(r' {2,}')
38 all = re_two_spaces.sub(" ", all)
39
40 re_comment = re.compile(r'<!--.+?-->')
41 all = re_comment.sub('', all)
42
43 re_nbsp = re.compile(r'(&nbsp;)+')
44 all = re_nbsp.sub(' ', all)
45
46 re_xmlns = re.compile(r' xmlns.+?>')
47 all = re_xmlns.sub('>', all)
48
49 re_blocktag = re.compile(r'<(p|li|div|table|tr|th|body|head|h.+?|title)>')
50 all = re_blocktag.sub("¥g<0¥n¥n", all)
51
52 re_blocktag2 = re.compile(r'<(body.+?|link.+?|br|head|div.+?|html.+?!DOCTYPE.+?)>')
53 all = re_blocktag2.sub("¥g<0¥n¥n", all)
54
55 re_blocktag3 = re.compile(r'<(meta.+?)>')
56 all = re_blocktag3.sub("¥g<0¥n", all)
57
58 re_if = re.compile(r'<!¥[(if.+?|endif)¥]>')
59 all = re_if.sub('', all)
60
61 re_div_style = re.compile(r'<div.+style=¥ layout-grid.+>')
62 all = re_div_style.sub('', all)
63
64 re_class = re.compile(r' (class|lang)=[^ >]+')
65 all = re_class.sub('', all)
66
67 re_html = re.compile(r'<html>')
68 all = re_html.sub('<html lang="ja">', all)
69
70 re_text = re.compile(r'(text¥-|mso¥-|margin¥-|font¥-family:|background:|tab¥-).+?¥')
71 all = re_text.sub("", all)
72
73 re_color_black = re.compile(r'(style=¥)?color:black.+?¥')
74 all = re_color_black.sub("", all)
75
76 re_span_langdir = re.compile(r'<span (lang|dir).+?>')
77 all = re_span_langdir.sub("<span>", all)
78
79 re_kigou1 = re.compile(r'¥&(sup2|0slash);')
80 all = re_kigou1.sub("", all)
81
82 re_msorm = re.compile(r' ¥!msorm')

```

henkan.py

```
83 all = re_msorm.sub("", all)
84
85 re_op = re.compile(r'</?o:p>')
86 all = re_op.sub("", all)
87
88 re_null_quotation = re.compile(r'[-a-z0-9]+=?¥ ¥')
89 all = re_null_quotation.sub('', all)
90
91 re_span_style = re.compile(r"<span style='font-[^]+?></span>")
92 all = re_span_style.sub('', all)
93
94 all = all.replace(' >', '>')
95
96 re_spannum = re.compile(r'</?span>')
97 all = re_spannum.sub(spannumber, all)
98
99 re_null_span = re.compile(r'<span(-?¥d+)>(.*?)</span¥1>')
100 m = re_null_span.search(all)
101 while m:
102     all = re_null_span.sub('¥g<2>', all)
103     m = re_null_span.search(all)
104 re_unspannum = re.compile(r'</?span>-?¥d+')
105 all = re_unspannum.sub('¥g<1>', all)
106
107 re_div = re.compile(r'</?div.*?>')
108 all = re_div.sub('', all)
109
110 re_karatag = re.compile(r'<([ ^>]+)[^>]* *¥1>')
111 m = re_karatag.search(all)
112 while m:
113     all = re_karatag.sub('', all)
114     m = re_karatag.search(all)
115
116 all = re_kaigyo.sub('¥n¥n', all)
117
118 all = '<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01 Transitional//EN">¥n' + all
119
120 all = all.encode('sjis')
121 print all
```