# Language Processing

Tuesday, July 26, 2022     4:12 PM

Assessing and verifying text and python functions with the characters found in this text
>    Whitespace

```
#Validate Whitespace behavior

1   white_space_chars = [k for k in char_counts if len(k.strip()) == 0]
2   print(white_space_chars)
                                                                          Python
[ ]
...   [' ', '\xa0', '\t']
```

>    Upper/Lower/IsAlpha/Punctuation

```
#Validate Upper/lower conclusion: using str.upper()/lower() functions is safe. No character in the dataset
causes an error when used in those functions, and the only characters that don't cooperate to a new case
are non-alphabetical characters such as numbers and punctuation. conclusion: using str.upper()==str.lower()
is a viable way to check if a character is alphabetical or not.

1   upperWITHlowerChars = set()
2   upper_chars = []
3   lower_chars = []
4   non_castable = []
5   error_casting = []
6   nons = []
7   found_as_one_case_only = []
8   for k in char_counts:
9 >     try: ...
22      except:
23          error_casting.append(k)
24  print('\nupper_chars:      ',sorted(upper_chars))
25  print('\nlower_chars:      ',sorted(lower_chars))
26  print('\nnon_castable:     ',sorted(non_castable))
27  print('\nerror_casting:    ',sorted(error_casting))
28  print('\nsilent fails:     ',sorted(nons))
29
30  print('\nupperWITHlowerChars:    ',sorted(upperWITHlowerChars))
31  for u,l in upperWITHlowerChars:
32      pair = [u,l]
33      unseen_possible_case = False
34 >    if l not in lower_chars and l not in non_castable: ...
37 >    if u not in upper_chars and u not in non_castable: ...
40 >    if unseen_possible_case: ...
42          # print("upper possible, but not present:      ",u)
43  print('\nfound_as_one_case_only:     ', sorted(found_as_one_case_only))
                                                                          Python
[ ]
```

```
upper_chars:      ['A', 'B', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'J', 'K', 'L', 'M', 'N', 'O', 'P',
'Q', 'R', 'S', 'T', 'U', 'V', 'W', 'X', 'Y', 'Z', 'Ñ', 'Ŋ', 'Ɓ', 'Đ', 'Ƴ']

lower_chars:      ['a', 'b', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'j', 'k', 'l', 'm', 'n', 'o', 'p',
'q', 'r', 's', 't', 'u', 'v', 'w', 'x', 'y', 'z', 'à', 'â', 'ç', 'è', 'é', 'ê', 'î', 'ī', 'ñ', 'ò',
'ô', 'ù', 'û', 'ŋ', 'y', 'ɓ', 'đ']

non_castable:     ['\t', ' ', '!', '"', '&', "'", '(', ')', '+', ',', '-', '.', '/', '0', '1', '2',
'3', '4', '5', '6', '7', '8', '9', ':', ';', '<', '=', '>', '?', '[', ']', '`', '\xa0']

error_casting:    []

silent fails:     []

upperWITHlowerChars:      [('A', 'a'), ('B', 'b'), ('C', 'c'), ('D', 'd'), ('E', 'e'), ('F', 'f'),
('G', 'g'), ('H', 'h'), ('I', 'i'), ('J', 'j'), ('K', 'k'), ('L', 'l'), ('M', 'm'), ('N', 'n'),
('O', 'o'), ('P', 'p'), ('Q', 'q'), ('R', 'r'), ('S', 's'), ('T', 't'), ('U', 'u'), ('V', 'v'),
('W', 'w'), ('X', 'x'), ('Y', 'y'), ('Z', 'z'), ('À', 'à'), ('Â', 'â'), ('Ç', 'ç'), ('È', 'è'),
('É', 'é'), ('Ê', 'ê'), ('Î', 'î'), ('Ï', 'ï'), ('Ñ', 'ñ'), ('Ò', 'ò'), ('Ô', 'ô'), ('Ù', 'ù'),
('Û', 'û'), ('Ŋ', 'ŋ'), ('Ɓ', 'ɓ'), ('Đ', 'đ'), ('Ƴ', 'y')]

found_as_one_case_only:       [(None, 'à'), (None, 'â'), (None, 'ç'), (None, 'è'), (None, 'é'),
(None, 'ê'), (None, 'î'), (None, 'ī'), (None, 'ò'), (None, 'ô'), (None, 'ù'), (None, 'û')]
```

```python
1  alpha_chars = [x for x in char_counts.keys() if x.upper() != x.lower()]
2  stralpha = [x for x in alpha_chars if x.isalpha()]
3  assert stralpha == alpha_chars, 'note that str.isalpha does NOT work safely here'
```
Python

Hex/Unicode values, and implications on regex [a-z] style ranges

```python
1  [(c,code_point(c)) for c in sorted(non_castable)]
2
3  # ('\t', 'U+0009' -> ('`', 'U+0060')
4  #  ('\xa0', 'U+00A0'))
5  # ('A', 'U+0041' -> 'û', 'U+00FB')
6  # ('Ŋ', 'U+014A' -> 'y', 'U+01B4')
7
```

Validating Regex basics

```python
#Validate Regex Behavior

1   impossible_char = '\u0008' #utf backspace (\u0008) is unlikely to appear in a docx, and did not
2   s = impossible_char.join(char_counts.keys())
3   re_results = [False]*len(char_counts.keys())
4   for i, k in enumerate(char_counts):
5       pattern = re.escape(k)
6       # print(s)
7       try:
8           m = re.search(pattern,s) #type: ignore
9           corrected_ind = m.start()/2
10          # print(corrected_ind)
11      except: print('exception: ',repr(i))
12      # print(corrected_ind)
13      if i == corrected_ind:
14          re_results[i] = True
15      else: print('failure: ',repr(i))
16  print(all(re_results))
```
Python

```
...    True
```