

C2_W1_lecture_nb_02_candidates_from_edits

December 13, 2023

1 NLP Course 2 Week 1 Lesson : Building The Model - Lecture Exercise 02

Estimated Time: 20 minutes # Candidates from String Edits Create a list of candidate strings by applying an edit operation ### Imports and Data

```
[ ]: # data
word = 'dearz' #
```

1.0.1 Splits

Find all the ways you can split a word into 2 parts !

```
[ ]: # splits with a loop
splits_a = []
for i in range(len(word)+1):
    splits_a.append([word[:i],word[i:]])

for i in splits_a:
    print(i)
```

```
[ ]: # same splits, done using a list comprehension
splits_b = [(word[:i], word[i:]) for i in range(len(word) + 1)]

for i in splits_b:
    print(i)
```

1.0.2 Delete Edit

Delete a letter from each string in the `splits` list. What this does is effectively delete each possible letter from the original word being edited.

```
[ ]: # deletes with a loop
splits = splits_a
deletes = []
```

```

print('word : ', word)
for L,R in splits:
    if R:
        print(L + R[1:], ' <-- delete ', R[0])

```

It's worth taking a closer look at how this is executing a 'delete'. Taking the first item from the splits list :

```

[ ]: # breaking it down
print('word : ', word)
one_split = splits[0]
print('first item from the splits list : ', one_split)
L = one_split[0]
R = one_split[1]
print('L : ', L)
print('R : ', R)
print('*** now implicit delete by excluding the leading letter ***')
print('L + R[1:] : ', L + R[1:], ' <-- delete ', R[0])

```

So the end result transforms 'dearz' to 'earz' by deleting the first character. And you use a **loop** (code block above) or a **list comprehension** (code block below) to do this for the entire splits list.

```

[ ]: # deletes with a list comprehension
splits = splits_a
deletes = [L + R[1:] for L, R in splits if R]

print(deletes)
print('*** which is the same as ***')
for i in deletes:
    print(i)

```

1.0.3 Ungraded Exercise

You now have a list of *candidate strings* created after performing a **delete** edit. Next step will be to filter this list for *candidate words* found in a vocabulary. Given the example vocab below, can you think of a way to create a list of candidate words ? Remember, you already have a list of candidate strings, some of which are certainly not actual words you might find in your vocabulary ! So from the above list **earz**, **darz**, **derz**, **deaz**, **dear**. You're really only interested in **dear**.

```

[ ]: vocab = ['dean', 'deer', 'dear', 'fries', 'and', 'coke']
      edits = list(deletes)

print('vocab : ', vocab)
print('edits : ', edits)

```

```
candidates=[]

### START CODE HERE ###
#candidates = ?? # hint: 'set.intersection'
### END CODE HERE ###

print('candidate words : ', candidates)
```

Expected Outcome:

```
vocab : ['dean', 'deer', 'dear', 'fries', 'and', 'coke'] edits : ['earz', 'darz', 'derz', 'deaz', 'dear']
candidate words : {'dear'}
```

1.0.4 Summary

You've unpacked an integral part of the assignment by breaking down **splits** and **edits**, specifically looking at **deletes** here. Implementation of the other edit types (insert, replace, switch) follows a similar methodology and should now feel somewhat familiar when you see them. This bit of the code isn't as intuitive as other sections, so well done! You should now feel confident facing some of the more technical parts of the assignment at the end of the week.