

AI9000

Alex Lang

Conor Rogers

Confusion

55%

We did not understand
Dep. Parsing OR Con. Parsing

```
# unzip and read story from zip file
def unzip_corpus(input_file, name):
    zip_archive = zipfile.ZipFile(input_file)
    contents = [zip_archive.open(fn, 'r').read().decode('utf-8')
                 for fn in zip_archive.namelist() if fn == name]
    return ''.join(contents)

# get bag of words
# Breaks sentences into set of tokenized sentences, removing stopwords
def get_bow(tagged_tokens, stopwords=""):
    if stopwords == "":
        return set([t[0].lower() for t in tagged_tokens])
    else:
        return set([t[0].lower() for t in tagged_tokens if t[0].lower() not in stopwords])

# The standard NLTK pipeline for POS tagging a document
def get_sentences(text):
    sentences = nltk.sent_tokenize(text)
    sentences = [nltk.word_tokenize(sent) for sent in sentences]
    sentences = [nltk.pos_tag(sent) for sent in sentences]

    return sentences

# Returns lemma of word added
def lemmatizer(tokens):
    lem_tokens = []
    # This little bit is because wordnet lemmas don't play nicely with things verb infinitives..... [very rough fix
    second_form_same_vinfinitive = [['felt', 'feel'], ['fell', 'fall'], ['stood', 'stand'], ['flattered', 'flatter'], ['flattery', 'flatter']]

    vinfinitive_check = [a for (a,b) in second_form_same_vinfinitive]
    for token in tokens:
        for (a,b) in second_form_same_vinfinitive:
            if a == token:
                lem_tokens += [b]
            if token not in vinfinitive_check:
                lem_tokens += [WordNetLemmatizer().lemmatize(token, 'v')]

    return lem_tokens

# qtokens: is a list of pos tagged question tokens with SW removed
# text: list of a list of pos tagged story sentences
# stopwords is a set of stopwords
# matches words to sentences for each text and returns the best answer
def baseline(qbow, text, stopwords):
    # collect all the candidate answers
    answers = []
    qbow = set([nltk.LancasterStemmer().stem(word) for word in qbow])
    qbow.update(set(lemmatizer(qbow)))
    print(qbow)
    for f in text:
        for sent in f:
            # A list of all the word tokens in the sentence
            sbow = get_bow(sent, stopwords)

            # stem all questions and sentences for better results
            sbow = set([nltk.LancasterStemmer().stem(word) for word in sbow])
            sbow.update(set(lemmatizer(sbow)))

            # and then add the other
            print(sbow)

            # Count the # of overlapping words between the Q and the A
            # & is the set intersection operator
            overlap = len(qbow & sbow)
            print(c.ORGREEN + "overlap: " + c.ENDC + str(overlap))

            answers.append((overlap, sent))

    # Sort the results by the first element of the tuple (i.e., the count)
    # Sort answers from smallest to largest by default, so reverse it
    answers = sorted(answers, key=operator.itemgetter(0), reverse=True)
    #print(answers)

    # Return the best answer
    best_answer = (answers[0])[1]

    return best_answer

# reads file and finds best sentence
def find_best_sentence(question, fnames):
    print(c.ORGREEN)
    print(fnames)

    # raw story / sch
    dataset = "hw5_dataset"
    text = [unzip_corpus(dataset + ".zip", dataset + "/" + f) for f in fnames]

    # get words for every sentence in sentence
    stopwords = set(nltk.corpus.stopwords.words("english"))
    #stopwords = ""
```

Rejection

```
# whois:
# r = r'^(S+/DT)?(S+/J)?\s?(S+/NN)+\s?(S+/NN)+$

# 1. looks at words in question and decides whether to look for words or POS
# WHO - looks for a POS NP "DT" "JJ" "NN"
search_words = [word for word, tag in q]
for 'who' in search_words:
    print('who')
    r = r'^(S+/DT)?\s?(S+/JJ)+\s?(S+/NN)+\s?(S+/NN)\s?(S+/NN)??'
    #default.

# WHAT - tricky, reads Q last 2 words, then searches for words after them
elif 'what' in search_words:
    if 'did' in search_words:
        #['that', 'IN'], ['the', 'DT'], ('lion', 'NN'), ('was', 'VBD'), ('friendly', 'RB')]
        #('a', 'DT'), ('large', 'JJ'), ('stainless', 'NN'), ('steel', 'NN'), ('bow!', 'NN'), ('of', 'IN'), ('water', 'NN')
        #('the', 'DT'), ('house', 'NN'), ('of', 'IN'), ('the', 'DT'), ('narrator', 'NN')
        # r = r'^(S+/NP)(S+/DT)?\s?(S+/JJ)+\s?(S+/NN)+\s?(S+/VBW)?\s?(S+/RB)?$'
        print('what did')
        r = r'^(S+/TO)?\s?(S+/OT)?\s?(S+/VBW)?\s?(S+/JJ)+\s?(S+/NNW)+\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/IN)\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/DT)?\s?(S+/NNW)?$'
        #r = r'^(S+/TO)?\s?(S+/OT)?\s?(S+/VBW)?\s?(S+/NN)+\s?(S+/JJ)+\s?(S+/NNW)?\s?(S+/DT)?\s?(S+/NNW)?\s?(S+/RBW)?$'
    elif 'happened' in search_words:
        #['police', 'NNS'], ('cars', 'NNS'), ('were', 'VBD'), ('burned', 'VBN')
        print('what happened')
        r = r'^(S+/OT)?\s?(S+/JJ)+\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/VBW)?\s?(S+/VBW)?\s?(S+/VBW)?$'
    elif 'was' in search_words:
        print('what was')
        #['police', 'NNS'], ('cars', 'NNS'), ('were', 'VBD'), ('burned', 'VBN')
        #['that', 'IN'], ('the', 'DT'), ('large', 'JJ'), ('and', 'CC'), ('large', 'JJ'), ('dish', 'NN')]
        r = r'^(S+/DT)?\s?(S+/JJ)+\s?(S+/CD)?\s?(S+/JJ)+\s?(S+/NNW)+\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/NNW)?\s?(S+/NNW)?$'
```

```
WHERE = looks for POS tag "IN" "DT" "NN" ?
elif 'where' in search_words:
    print('where')
    # ('in', 'IN'), ('a', 'DT'), ('flat', 'JJ'), ('and', 'CC'), ('large', 'JJ'), ('dish', 'NN')
    r = r'(\S+/IN)+\s?(\S+/DT)\s?(\S+/JJ)\s?(\S+/CC)\s?(\S+/JJ)+\s?(\S+/NNW?)+

# WHEN - tricky as well, looks for POS "IN"
elif 'when' in search_words:
    print('when')
    # ('a', 'DT'), ('few', 'JJ'), ('years', 'NNS'), ('ago', 'RB')
    r = r'(\S+/DT)\s?(\S+/JJ)+\s?(\S+/NNW?)+\s?(\S+/RB)?

# WHY - looks for the word "because" or just the Q and words after it
elif 'why' in search_words:
    print('why')
    # ('in', 'IN'), ('order', 'NN'), ('for', 'IN'), ('the', 'DT'), ('birds', 'NNS'), ('to', 'TO'), ('wait', 'VB')
    # ('because', 'IN'), ('it', 'PRP'), ('owned', 'VBD'), ('a', 'DT'), ('group', 'NN'), ('of', 'IN'), ('trees', 'NNS'), ('and', 'CC'), ('near', 'IN'), ('the', 'DT'), ('house', 'NN')
    r = r'(\S+/IN)+\s?(\S+/NNW?)\s?(\S+/VBD)\s?(\S+/DT)\s?(\S+/NNW?)\s?(\S+/IN)+\s?(\S+/RB)\s?(\S+/JJ)\s?(\S+/NNW?)+\s?(\S+/TO)\s?(\S+/VBW?)\s?(\S+/DT)\s?(\S+/NNW?)\s?(\S+/DT)\s?(\S+/NNW?)\s?(\S+/DT)\s?(\S+/NNW?)+
    # r = r'(\S+/because)+
    # looks for because
```

```
elif 'how' in search_words:
    print('how')
    r = r'(\$/wowthisisareallylongregexexpressionthatworksgreatanditisawesome)*\$?(\$/RB)*\$?'
```

```
answers = get_phrase(s_proc_nostem, r)
```

```
for answer in answers:
    remove = False
    for word in answer:
        q_senty = [word.lower() for word, tag in q]
        q_senty = [w for w in q_senty if w not in stopwords]
        if word in q_senty:
            remove = True
    if remove:
        answers.remove(answer)
```

```
# print(answers)

if len(answers) > 0:
    answer = ' '.join(w for w in answers[0])
else:
    answer = ''

# if s_matches != []:
#     answer = s_matches[high[0]]

print(c.OKGREEN + "ANSWER: " + c.ENDC, answer)
```

Rejection

58%

We used REGEX
& a Dumb Luck Algorithm

Acceptance

65%

We switched to
Constituency Parsing

```
import argparse, re, nltk, math, sys
from nltk.tree import Tree
import chunky
import process_questions

# Read the constituency parse from the line and construct the Tree
def read_con_parsers(parfile):
    fh = open(parfile, 'r')
    lines = fh.readlines()
    fh.close()
    return [Tree.fromstring(line) for line in lines]

# See if our pattern matches the current root of the tree
def matches(pattern, root):
    # Base cases to exit our recursion
    # If both nodes are null we've matched everything so far
    if root is None and pattern is None:
        return root

    # We've matched everything in the pattern we're supposed to (we can ignore the extra
    # nodes in the main tree for now)
    elif pattern is None:
        return root

    # We still have something in our pattern, but there's nothing to match in the tree
    elif root is None:
        return None

    # A node in a tree can either be a string (if it is a leaf) or node
    plabel = pattern if isinstance(pattern, str) else pattern.label()
    rlabel = root if isinstance(root, str) else root.label()

    # If our pattern label is the * then match no matter what
    if plabel == "*":
        return root
    # Otherwise they labels need to match
    elif plabel == rlabel:
        # If there is a match we need to check that all the children match
        # Minor bug (what happens if the pattern has more children than the tree)
        for pchild, rchild in zip(pattern, root):
            match = matches(pchild, rchild)
            if match is None:
                return None
        return root
    else:
        return None

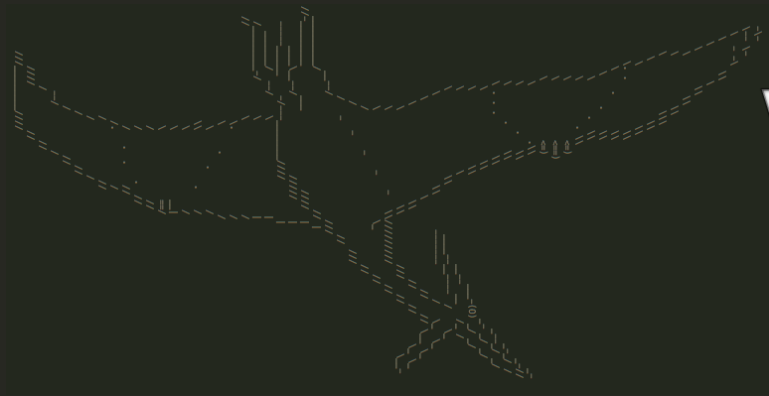
def pattern_matcher(pattern, tree):
    for subtree in tree.subtrees():
        node = matches(pattern, subtree)
        if node is not None:
            return node
    return None

def subtree_master(pattern, tree):
    subtree = pattern_matcher(pattern, tree)

    if subtree is None:
        # pattern = nltk.ParentedTree.fromstring("(NP (w) )")
        # subtree = pattern_matcher(pattern, tree)
        # print(subtree)
        # print(" ".join(subtree.leaves()))
        return 'nopes'

    print(subtree)
    print(" ".join(subtree.leaves()))
    return subtree

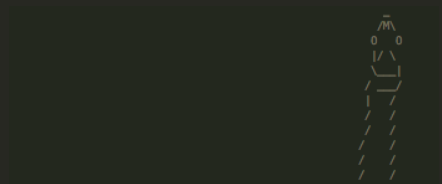
def q_determine(question, tree):
    search_words = question.lower().split()
    # print(search_words)
    if 'who' in search_words:
        print('who')
```



Victory



68%



Held-Out Data

Recall

64.4%

Precision

63.2%

F-score

63.8%

Reflections

Q?

We don't 'know' what the Q is asking
We treat each 'wh' Q type accordingly

Observations

Recall

What Q is asking for &
Finding Best Sentence

Precision

Answering Q