

In [61]: FINAL_PUNC_CLEANER = str.maketrans(dict([(c, '') for c in '!"#\$%^&'()*+,-./:;<=>?[\]^_`{|}~--@']])

```
def clean_text(text):
    # odd chars
    # will generate more ' so do it first
    text = re.sub('\'', '', text) # ?
    text = re.sub('\"', '', text) # special single quote
    text = re.sub('\'', '', text) # special single quote
    text = re.sub('\"', '', text) # special double quote
    text = re.sub('\'', '', text)
    text = re.sub('\"', '', text)
    text = re.sub('\'', '', text)

    # shortcuts
    text = re.sub('\s', ' ', text)
    text = re.sub('whats', ' what is ', text)
    text = re.sub('ve', ' have ', text)
    text = re.sub('can\'', 'can not', text)
    # this one is tricky do it in order
    text = re.sub('wouldn\'', 'would not', text)
    text = re.sub('n\'', ' not ', text)
    text = re.sub('i\'m', 'i am', text)
    text = re.sub('re', ' are ', text)
    text = re.sub('d', ' would ', text)
    text = re.sub('ll', ' will ', text)
    text = re.sub('e.g.', ' eg ', text)
    text = re.sub('b.g.', ' bg ', text)
    text = re.sub('e-mail', ' email ', text)
    text = re.sub('\s\)', ' ', text)

    # Numbers and measures are a true mess
    # 12,000 -> 12000
    text = re.sub('(?!=[0-9])\,(?!=[0-9])', '', text)

    # Quora is very used in India so rouple (rs) is often present
    text = re.sub("(?!=[0-9])rs ", " rs ", text)
    text = re.sub("rs(?!=[0-9])", " rs ", text)

    # stolen at kaggle : https://www.kaggle.com/currie32/the-importance-of-cleaning-text

# text = re.sub('[-FC-FJ]\s|/ ', ' disk ', text)
# text = re.sub('([d+](K))', ' [g<1000 ', text)
# very weird !!! these ones decrease the hit % WTF ?

# text = re.sub(r"(the[ls]+the[ls]+)?us(a)? ", " usa ", text)
# text = re.sub('([the[ls]+the[ls]+)?united state(s)?', ' usa ', text)

text = re.sub(r" uk ", " england ", text)
text = re.sub(r" improvement ", " improvement ", text)
text = re.sub(r" initially ", " initially ", text)
text = re.sub(r" dms ", " direct messages ", text)
text = re.sub(r" demonitization ", " demonetization ", text)
text = re.sub(r" activated ", " active ", text)
text = re.sub(r" kms ", " kilometers ", text)
text = re.sub(r" cs ", " computer science ", text)
text = re.sub(r" upvote ", " up vote ", text)
text = re.sub(r" iphone ", " phone ", text)
text = re.sub(r" \0rs ", " rs ", text)
text = re.sub(r" calendar ", " calendar ", text)
text = re.sub(r" ios ", " operating system ", text)
text = re.sub(r" programming ", " programming ", text)
text = re.sub(r" bestfriend ", " best friend ", text)
text = re.sub(r" iii ", " 3 ", text)
text = re.sub(r" banglore ", " bangalore ", text)
text = re.sub(r" j k ", " jk ", text)
text = re.sub(r" J\K\.", " jk ", text)

# some others
text = re.sub(r"60k", " 60000 ", text)
text = re.sub(r" e g ", " eg ", text)
text = re.sub(r" b g ", " bg ", text)
text = re.sub(r"\0s", "0", text)
text = re.sub(r" 9 11 ", "911", text)
text = re.sub(r"s(2)", " ", text)
text = re.sub(r" usa ", " America ", text)
text = re.sub(r" u s ", " America ", text)
text = re.sub(r"m ", " am ", text)

# will blank any of !"#$%&'()*+,-./:;<=>?[\]^_`{|}~--@
# Note the @ is dubious : won't we loose some emails
text = text.translate(FINAL_PUNC_CLEANER)

return text
```

```
small_train = new_preprocess(small_train)
#display(small_train[['common_words', 'new_common_words', 'uncommon_words_question1', 'new_uncommon_words_q
uestion1']].head(10))
sniff_changes(small_train)
```

Compute all features in one shot

Extract common words between question1 & question2

Extract uncommon words in question1

Extract Nb common words between question1 & question2

Extract Nb words in question1 not in common words

New common 118.768 % New uncommon 82.289 %

In [64]: # I do special stuff with \$ and rouple char
FINAL_PUNC_CLEANER = str.maketrans(dict([(c, '') for c in '!"#\$%^&'()*+,-./:;<=>?[\]^_`{|}~--@']])

```
def clean_text(text):
    # odd chars
    # will generate more ' so do it first
    text = re.sub('\'', '', text) # ?
    text = re.sub('\"', '', text) # special single quote
    text = re.sub('\'', '', text) # special single quote
    text = re.sub('\"', '', text) # special double quote
    text = re.sub('\'', '', text)
    text = re.sub('\"', '', text)
    text = re.sub('\'', '', text)

    # shortcuts
    text = re.sub('\s', ' ', text)
    text = re.sub('whats', ' what is ', text)
    text = re.sub('ve', ' have ', text)
    text = re.sub('can\'', 'can not', text)
    # this one is tricky do it in order
    text = re.sub('wouldn\'', 'would not', text)
    text = re.sub('n\'', ' not ', text)
    text = re.sub('i\'m', 'i am', text)
    text = re.sub('re', ' are ', text)
    text = re.sub('d', ' would ', text)
    text = re.sub('ll', ' will ', text)
    text = re.sub('e.g.', ' eg ', text)
    text = re.sub('b.g.', ' bg ', text)
    text = re.sub('e-mail', ' email ', text)
    text = re.sub('\s\)', ' ', text)

    # Numbers and measures are a true mess
    # 12,000 -> 12000
    text = re.sub('(?!=[0-9])\,(?!=[0-9])', '', text)

    # Quora is very used in India so rouple (rs) is often present
    text = re.sub("(?!=[0-9])rs ", " rs ", text)
    text = re.sub("rs(?!=[0-9])", " rs ", text)

    # stolen at kaggle : https://www.kaggle.com/currie32/the-importance-of-cleaning-text

# text = re.sub('[-FC-FJ]\s|/ ', ' disk ', text)
# text = re.sub('([d+](K))', ' [g<1000 ', text)
# very weird !!! these ones decrease the hit % WTF ?

# text = re.sub(r"(the[ls]+the[ls]+)?us(a)? ", " usa ", text)
# text = re.sub('([the[ls]+the[ls]+)?united state(s)?', ' usa ', text)

text = re.sub(r" uk ", " england ", text)
text = re.sub(r" improvement ", " improvement ", text)
text = re.sub(r" initially ", " initially ", text)
text = re.sub(r" dms ", " direct messages ", text)
text = re.sub(r" demonitization ", " demonetization ", text)
text = re.sub(r" activated ", " active ", text)
text = re.sub(r" kms ", " kilometers ", text)
text = re.sub(r" cs ", " computer science ", text)
text = re.sub(r" upvote ", " up vote ", text)
text = re.sub(r" iphone ", " phone ", text)
text = re.sub(r" \0rs ", " rs ", text)
text = re.sub(r" calendar ", " calendar ", text)
text = re.sub(r" ios ", " operating system ", text)
text = re.sub(r" programming ", " programming ", text)
text = re.sub(r" bestfriend ", " best friend ", text)
text = re.sub(r" iii ", " 3 ", text)
text = re.sub(r" banglore ", " bangalore ", text)
text = re.sub(r" j k ", " jk ", text)
text = re.sub(r" J\K\.", " jk ", text)

# some others
text = re.sub(r"60k", " 60000 ", text)
text = re.sub(r" e g ", " eg ", text)
text = re.sub(r" b g ", " bg ", text)
text = re.sub(r"\0s", "0", text)
text = re.sub(r" 9 11 ", "911", text)
text = re.sub(r"s(2)", " ", text)
text = re.sub(r" usa ", " America ", text)
text = re.sub(r" u s ", " America ", text)
text = re.sub(r"m ", " am ", text)

# units
text = re.sub(r"(\d+)kgs ", "Lambda m: m.group(1) + ' kg ', text) # e.g. 4kgs => 4 kg
text = re.sub(r"(\d+)k ", "Lambda m: m.group(1) + ' kg ', text) # e.g. 4kg => 4 kg
text = re.sub(r"(\d+)k ", "Lambda m: m.group(1) + '000 ', text) # e.g. 4k => 4000
text = re.sub(r"(\d+)\$", "Lambda m: m.group(1) + ' dollar ', text)
# This one is important in 2017
text = re.sub("donald trump", ' trump ', text)
text = re.sub('dollars', ' dollar ', text)
text = re.sub('quor', ' quora ', text)
text = re.sub("googling", " google ", text)
text = re.sub("googled", " google ", text)
text = re.sub("googleable", " google ", text)
text = re.sub("googles", " google ", text)

text = re.sub("rs ", " rs ", text) # []
text = re.sub("rs ", " dollar ", text)

# will blank any of !"#$%&'()*+,-./:;<=>?[\]^_`{|}~--@
# Note the @ is dubious : won't we loose some emails
# and $ is replaced by dollar before
text = text.translate(FINAL_PUNC_CLEANER)

return text
```

```
small_train = new_preprocess(small_train)
#display(small_train[['common_words', 'new_common_words', 'uncommon_words_question1', 'new_uncommon_words_q
uestion1']].head(10))
sniff_changes(small_train)
```

Compute all features in one shot

Extract common words between question1 & question2

Extract uncommon words in question1

Extract Nb common words between question1 & question2

Extract Nb words in question1 not in common words

New common 118.681 % New uncommon 82.213 %

Question that may have impact next steps How many unicode chars are present ?

In [65]: # Ugly but easy
def sniff_unicode(s):
 if re.sub('[\x00-\x7f]+', '', s) != s:
 return 1
 else:
 return 0

nb_unicode_train = train_dataframe['question1'].progress_apply(sniff_unicode).sum()/len(train_dataframe)
nb_unicode_train == train_dataframe['question2'].progress_apply(sniff_unicode).sum()/len(train_dataframe)
nb_unicode_challenge = challenge_dataframe['question1'].progress_apply(sniff_unicode).sum()/len(challenge_dataframe)
nb_unicode_challenge == challenge_dataframe['question2'].progress_apply(sniff_unicode).sum()/len(challenge_dataframe)
print_warning('Unicode train: %.3f challenge: %.3f' %(nb_unicode_train, nb_unicode_challenge))

Unicode train: 0.022 challenge: 0.021

Humpf, 0.02% of questions has some words fully in Unicode:

- awfully small.
- we are pushing the limits of this 'nb common words' feature and every little bit of information matters for the kaggle result
- Not clear if it does worth it ...
- Not clear also if we just replace it with nothing or with a generic placeholder like 'unicode_text' ...

I decide to not go so far (and a quick test showed kaggle score was slightly decreasing)