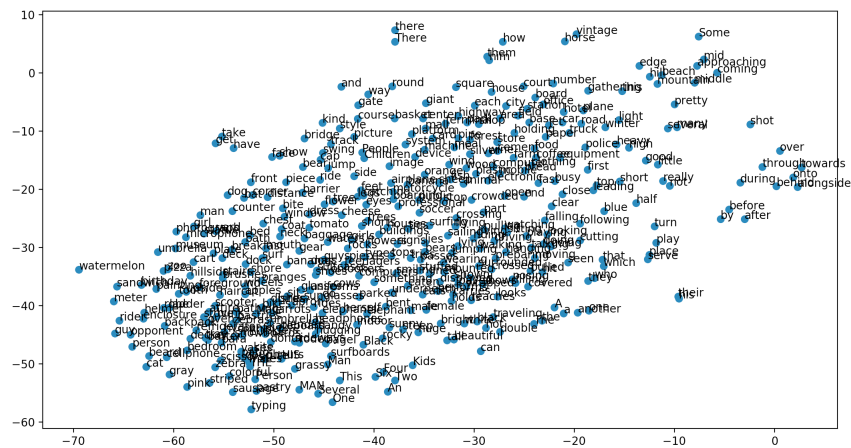
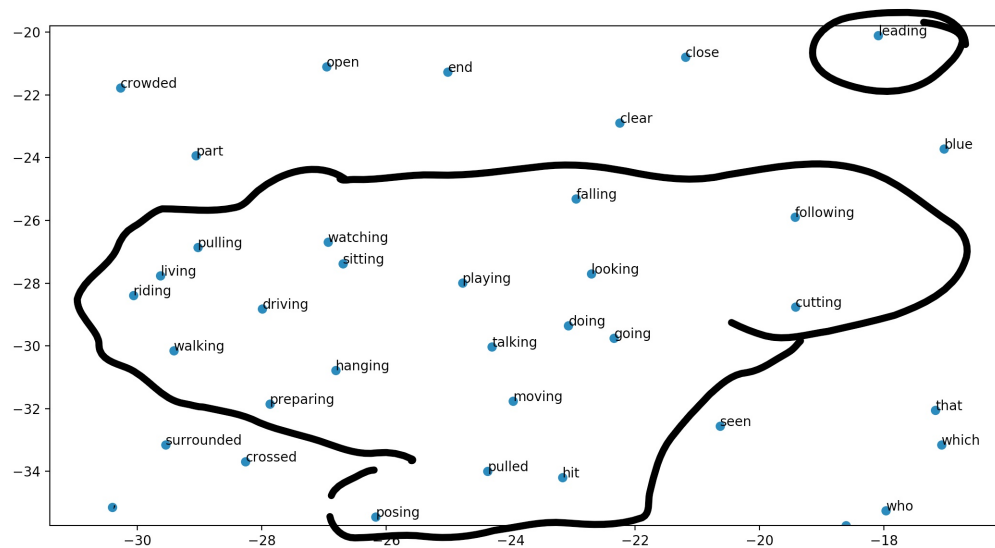


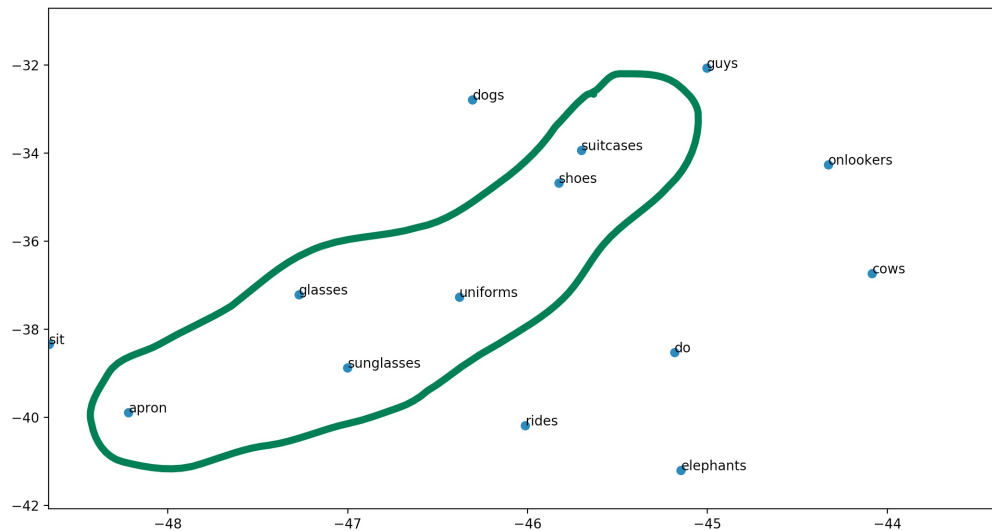
### Problem1 Visualize Word Embedding

I increase the amount of words from the original, I choose a subset of [200, 600] from the whole vocabulary.



For this part, I find verb with ing are grouped together as well as all nouns for apparel.





I select some word pairs from two groups, and calculate their cosine similarities in the 2-D diagram:

```
def similarity(x, y):
    return numpy.dot(x, y)/(numpy.linalg.norm(x) * numpy.linalg.norm(y))

#Verbs+ing group
talking = numpy.array([-24.29, -30.09])
doing = numpy.array([-23.07, -29.41])
going = numpy.array([-22.35, -29.85])
looking = numpy.array([-22.71, -22.77])

#apparel group
glasses = numpy.array([-47.24, -37.37])
uniforms = numpy.array([-46.37, -37.42])
shoes = numpy.array([-45.82, -34.72])

similarity(talking, doing)      -- 0.999902
similarity(doin, going)        -- 0.999747
similarity(doin, looking)      -- 0.992939
similarity(doin, glasses)      -- 0.972198
similarity(glasses, uniforms)  -- 0.999952
similarity(doin, uniforms)     -- 0.974427
```

Besides, I find that they also show cosine similarity in full dimension.

'Far' group means cosine similarity  $< 0.25$ , and 'Close' group means large cosine similarity  $> 0.75$ .

As shown in the following diagram, I calculate the cosine similarity to words 'doing'. I find that words has larger cosine similarity to 'doing' has stronger syntactic and meaning relationships(verb + ing form are more likely to be in close group). On the contrary, smaller cosine similarity means less relation close.

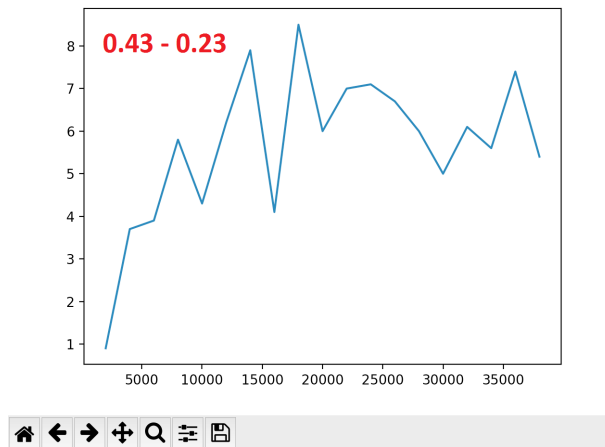
```
shimengdideMacBook-Pro:~$ cat a2-release mengdishi
Far-----
middle
this
ride
cap
cat
decker
donut
mall
male
Children
Close-----
pulling
something
jumping
playing
displayed
cutting
going
filling
good
doing
!
```

For the noun 'uniforms', the conclusion are the same. As we can see, the noun 'apparel' words are in close groups. In the 'Far' group, many verb words(including verb + ing).

```
Far-----
middle
this
ride
cap
decker
donut
crossed
can
following
pulling
Close-----
car
pieces
cows
airplane
clothing
suitcases
buildings
shoes
tops
rocks
shimengdideMacBook-Pro:a2-release mengdishi$
```

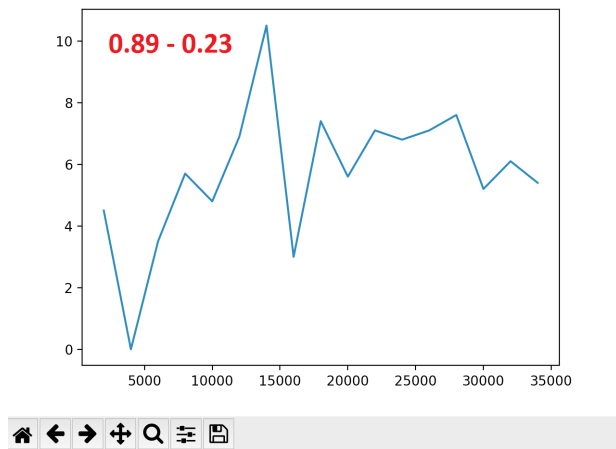
## Problem2 Train and Evaluate Model

1. My 4 setting of parameters of learning rate and momentum are:  
[0.43, 0.23]  $\rightarrow$  bleu score = 8.5

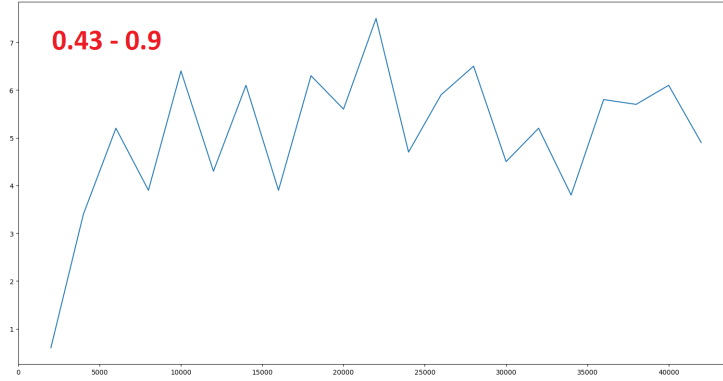


[0.89, 0.23]  $\rightarrow$  bleu score = 10.5

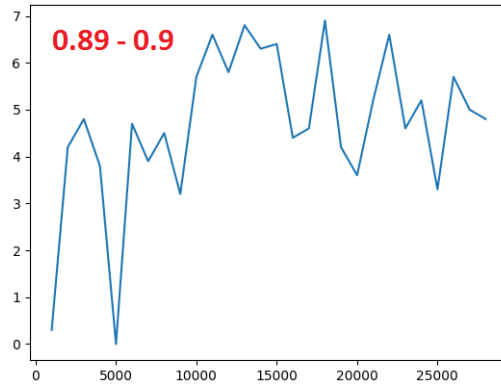
---



$[0.43, 0.9] \rightarrow \text{bleu score} = 7.5$



$[0.89, 0.9] \rightarrow \text{bleu score} = 6.9$



From the experiments, I find that large momentum(0.9) will cause diagram oscillate vertically. Especially when learning rate and momentum are large at the same time, it will cause each time gradient update change a lot, hence more likely to be affected by noise, which could be reflected by the 4th diagram, whose oscillation and amplitude change obviously. To the learning rate, smaller learning rate smooth curve, but more likely to be trapped in local minimum or maybe saddle point. On the contrary, large learning rate may cause it update more quickly, more likely to avoid local minimum. Therefore, it make sense that the 2nd setting of parameters  $[0.89, 0.23]$  get the highest bleu scores.

2. From result above, I set learning rate as 0.89 and momentum as 0.23 as my best model. And get the bleu score of 9.8 in testing:

```
(228, u'A group of people .')
(229, u'A group of people .')
(230, u'A man in a snow .')
(231, u'A group of people .')
(232, u'A man of a man in the air .')
(233, u'A man standing in the air .')
(234, u'A man of a man .')
(235, u'A group of people standing in the air .')
(236, u'A couple of people .')
(237, u'A man of people standing in the air .')
(238, u'A man of a man in the air .')
(239, u'A group of people .')
(240, u'A man of people .')
(241, u'A group of people .')
(242, u'A man is sitting at a table .')
(243, u'A group of people .')
(244, u'A man of a man .')
(245, u'A man in a man with a unk .')
(246, u'A group of people .')
(247, u'A man of a man .')
(248, u'A group of people .')
(249, u'A group of people .')
BLEU score = 9.8
shimengdideMacBook-Pro:a2-release mengdishi$
```

3. By setting beam width as 2, I get the following output:  
It actually keep computing conditional probability given front positions candidates. As we can see can see from the picture, for the 3rd words, although to the 3rd word, 'man' has larger probability follow 'A' than 'woman' follow 'A', so 'A man' has a bigger chance to be the first 2 words. However, for the 3rd word, although 'is' has bigger probability than 'standing', it is actually not the 3rd word in the final interpreter. Since, the probability of a interpreter sentence use a chain rule to calculate, which means it will consider every words probability to get the total probility for the sentence. Hence, althouh some words has tight bound locally, they may not get together at the final.

```

Packing up...
top 2 (word, log probability) @ time step 1:
(A, -7.97537552771)
(a, -8.81074235544)
top 2 (word, log probability) @ time step 2:
(man, -9.54097215498)
(woman, -11.5202657635)
top 2 (word, log probability) @ time step 2:
(man, -8.15602417347)
(woman, -10.6969137621)
top 2 (word, log probability) @ time step 3:
(is, -7.93423448455)
(standing, -10.9998323211)
top 2 (word, log probability) @ time step 3:
(standing, -7.87199485525)
(sitting, -11.0314520219)
top 2 (word, log probability) @ time step 4:
(on, -8.28807681355)
(sitting, -9.82658793846)
top 2 (word, log probability) @ time step 4:
(on, -9.0503801185)
(in, -10.3620144983)
top 2 (word, log probability) @ time step 5:
(a, -8.40483376602)
(the, -8.08399286414)
top 2 (word, log probability) @ time step 5:
(a, -8.07853637516)
(the, -8.72391904876)
top 2 (word, log probability) @ time step 6:
(tennis, -7.32941078452)
(unk, -12.732499452)
top 2 (word, log probability) @ time step 6:
(game, -7.29955852484)
(ball, -13.0675382862)
top 2 (word, log probability) @ time step 7:
(., -7.89241439625)
(court, -13.8109018341)
top 2 (word, log probability) @ time step 7:
(., -7.26173839020)
(court, -13.1440667436)
top 2 (word, log probability) @ time step 8:
(<end>, -14.832272273)
(unk, -17.8771426785)
top 2 (word, log probability) @ time step 8:
(<end>, -14.4727508194)

```

- (4) Just as shown at the above example, some words may have very tight bound, and hence will have very high probability with each other, which will cause bad captions if we just choose sentences based on information locally. But by using beam search, the probability of a word given several front words, which will gain more meanings to whole sentences, which will have more chance to make good caption.