

Word Embedding by Using Word2Vec and BERT

Github link: https://github.com/Allen-ZKW/NLP_HZAU/tree/task7

Abstract

In this task, we use two algorithms to complete Word Embedding. These two algorithms all try to transform words to vectors, these vectors will reflect the closeness and similarity between two different words.

Principle

Both of these two algorithms' target is translating word to type of information which can be understood by computers and programs. We can use 'translate data' to calculate distance between two different words. BERT algorithm takes word in different context into concern, so the result of BERT will include more information. This model is trained by forwards information and afterwards information.

Measure

Data Preparation

```
unzip data/litcovid-trainingdata.zip
```

Environmental Preparation

```
conda create -n NLP_task7 python=3.8
conda activate NLP_task7
conda install pytorch torchvision torchaudio cpuonly -c pytorch
conda install --yes --file requirements.txt
```

Word2Vec

```
python3 skip_gram_basic.py
```

BERT

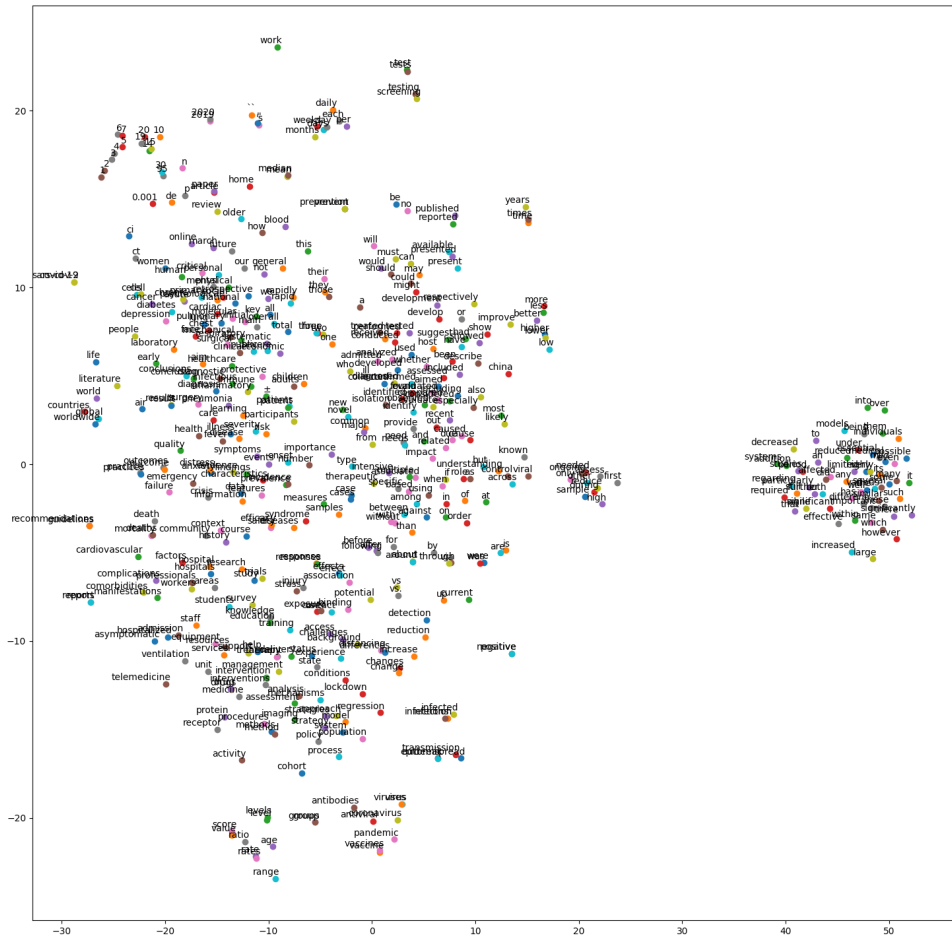
```
python3 Bert_4_Litcovid_wordEmbedding.py
```

Result

BERT Vectors

the 0.35628464818000793 -0.5250380635261536 -0.2492087334394455 -0.8056405186653137 -0.9308571815490723 0.16308681666851044 0.13157081604003906 -0.0601432062685489 -0.08181028813123703 0.5495608448982239 0.3162820637226105 0.044546108692884445 0.7944148182868958 0.3056371808052063 -0.31843802332878113 -0.1738629937171936 0.156353; 47095 -0.4021967947483063 1.411456823348999 -0.40982916951179504 0.3186410367488861 -1.1165008544921875 -0.05547158420085907 0.7745139002799988 0.43387165665626526 -0.2; .2594403326511383 -1.0255253314971924 0.7162691950798035 -0.3134252727031708 -0.33627140522003174 0.20484042167663574 -0.7187500596046448 -0.9545416235923767 -1.2280001; 1.018480777404785 0.5656768083572388 0.29735568165779114 -0.6946800947189331 -1.0299690961837769 1.2491761445999146 0.6057869791964558 0.5968462824821472 -0.032061982; 7 -1.5654470920562744 0.10546226799488068 -0.3829757869243622 -1.0525271892547607 0.10693001747131348 0.7864586114883423 0.5608898997306824 0.24788601696491241 0.526809 864502906799 -0.3442825973033905 0.27733975648880005 -0.7152450084666279 -0.24689675867557526 -0.7228854298591614 -1.2055928707122803 0.0033024270087480545 1.121882319; 32605 -1.6373316049575806 1.6440918445587158 -0.43068304657936096 -0.473300576210022 -1.2751322984695435 -0.25246310234069824 1.0536020994186401 -0.9760722517967224 -0.0 2385823726654053 -0.4132073223590851 0.01643933542072773 0.216070756316185 0.1842733770608902 0.47540250420570374 -0.5665438771247864 -0.2554280757904053 1.06030929088 62297058105 0.017322789877653122 -0.1162325198888779 -0.7282307147979736 -0.9285045862197876 -0.8524077534675598 0.11137823760509491 1.0579420328140259 0.408930689096 97026824951 -0.03119588829576969 0.10418453812599182 0.5999304056167603 -1.0119234323501587 -0.3879777193069458 -0.2799862325191498 -1.423446774482727 0.10861497372388; -0.2590559422969818 0.5961945652961731 0.47084879875183105 0.4812585413455963 -1.0473358631134033 0.00704399636015296 0.5069509744644165 0.07217913866043091 0.44964727 28891 -0.5169460773468018 -1.1743741035461426 -0.7117654651069641 -0.47802042961120605 -1.6063034534454346 -0.16897273063659668 1.5204949378967285 -0.1661432683467865 -0 80347061157 0.4028494656085968 -0.18535855412483215 0.7548332810401917 0.31232325705528259 0.6593093872070312 0.3685827851295471 0.21299351751804352 0.9169759750366211 04139697551727 -0.29513370990753174 0.5013274550437927 -0.21574245393276215 -0.5136917233467102 -0.9796431064605713 0.28284478187561035 0.15776784718036652 0.992968142; of 0.2614080011844635 -0.05700542405247668 0.050627123564481735 -0.12816961109638214 -0.5257653594017029 0.06948389858007431 0.34602680802345276 0.34332260489463; 928033351898 -0.11494740843772888 0.4615645110607147 0.34988754987716675 0.1981586515903473 0.06026952341198921 -0.8913010954856873 0.10875722020864487 0.0937142521142 095 -0.10814546048641205 -0.294975221157074 0.48537203669548035 -0.7861011624336243 0.3077158033847809 -0.875210702419281 0.5334144234657288 0.8603790998458862 0.183168. 0.014857389964163303 -0.19131046533584595 -0.14025750756263733 0.22238968312740326 -0.11352309584617615 0.12927047908306122 -0.2540449798107147 0.40016505122184753 -0.76 0.03420163318514824 -0.4061053991371749 -0.037243809551000595 0.3396824598312378 0.5161757469177246 -0.4832783043384552 -0.5008599162101746 0.540671706199646 0.06227. -0.5704774260520935 -0.10188032686710358 -0.5763516426086426 -0.8384292721748352 -0.2848418653011322 0.058282624930143356 -0.0714394822716713 0.31795915961265564 -0.126; 42 0.23032259941101074 0.2935047149658203 0.11147939413785934 0.06336495280265808 0.18873421847820282 0.16351547837257385 -0.36865437030792236 -0.1314295083284378 -0.56 -0.7546107769012451 0.06165650486946106 0.47786179184913635 -0.08768445998430252 0.32733067870140076 0.9853317141532898 -0.045827437192201614 -0.12225066870450974 -0.53 74774 0.058116666972637177 -0.0871891975402832 -0.4840514063835144 -0.21969972550868988 -0.00875602662563324 -0.026952695101499557 0.24102148413658142 0.02383398078382; 4354038239 0.11455033719539642 -0.20366010069847107 -0.18621009588241577 -0.28095537424087524 -0.07885187864303589 0.11217619478702545 -0.1198195070028305 -0.457789272; 01382040977478 0.0042733740992844105 -0.5538504123687744 0.4073474705219269 0.33703216910362244 -0.38002899289131165 -0.6938464045524597 -0.011663113720715046 -0.09683; 3527934551239 -0.30656880140304565 -0.05122643709182739 0.46377986669540405 0.1711607426404953 0.3399849832057953 0.4210023581981659 -0.7595409750938416 -0.04228069633 07845306396 0.32652902603149414 -0.1293913573026657 0.46423065662384033 0.29337844252586365 -0.24779483675956726 -0.033143870532512665 0.15377464890480042 0.3256935179 172455 0.044297799468040466 0.45039498805999756 -0.2799057066440582 -0.230113700032342 0.16503505408763885 0.11548270285129547 -0.0797380805015564 0.4481073021888733 - 4158 -0.3439595401287079 -0.14692458510398865 0.45710164308547974 -0.19501779973506927 0.30751246213912964 -0.14463791251182556 -0.25958913564682007 -0.2648999691009521 2010193 -0.3682197034358978 0.29388511180877686 -0.007812407799065113

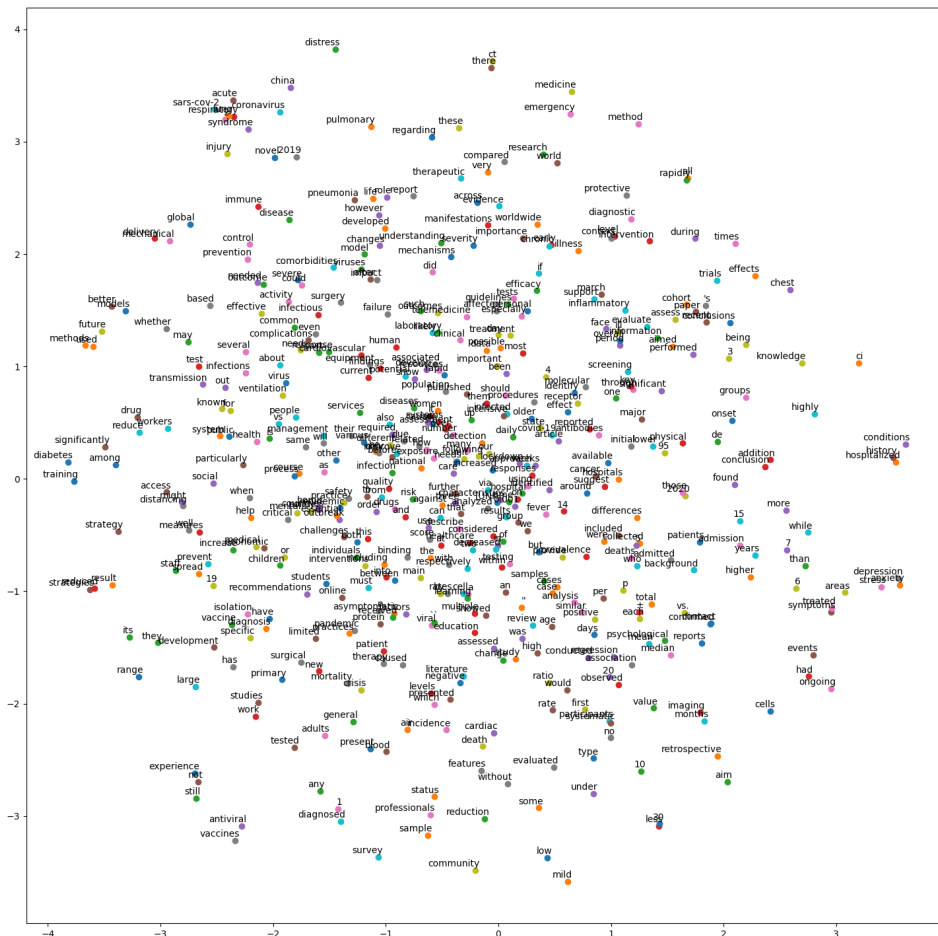
BERT Visualization



Word2Vec Vectors

[UNK] -0.04881330579519272 -0.0379069391010898 -0.05652253329753876 -0.09073057370729065 -0.11317484825849533 0.05839929729700885 0.020258576345443726 0.04351604357
14718006134 -0.011797694489359856 0.10310874134302139 0.01187390339958660 0.11605213582515717 -0.00109352682209015 -0.04330721125006676 0.03871
611858368
the -0.143330588968680573 -0.1587565541267395 0.14328989847507477 -0.3877052068710327 -0.10853473864144012 -0.03758662149310110 0.04662489991052246 0.02643380529
78235 -0.04327136278152466 -0.07235678285360336 -0.3090354800224304 0.12058231979608536 -0.02043553069238943 0.10384968668222427 0.28311672806739807 0.135265320539474
966949463
and 0.03092149654090103 0.06394444406032562 0.22052952647209167 0.12312621623277664 0.13519074022769928 -0.33134761452674866 -0.08193377405405045 0.17148441076278
10996109992265701 0.08987542241811752 -0.06514536589384079 0.0856313705444336 -0.06625906378030771 0.11863439530134201 -0.1884736567735672 -0.07212437689304352 0.09053
7
-0.022927175748827267 -0.27186140418052673 -0.1382111309588623 0.2388809472322464 -0.0699752010450636 0.0157379353594894 -0.2610098421576369 -0.03190140140
710476875305 0.1495684082764587 -0.26795533299446106 -0.29544875025749207 0.1886664034524515 -0.15007612109184265 0.13896331232738495 -0.015536155551671982 0.1629645
803672458976507
in 0.07111769165969016 -0.09316769242286682 -0.024860944598913193 0.12261395901441574 -0.32726961374282837 0.15798386931419373 -0.0223913059949875 -0.03950957953
1018983211368322232 -0.1676515047818298 -0.1708147227641296 0.033556707230388214 -0.2302724027633667 0.1060002226527037 0.17270736639312714 0.2785123884677887 -0.0613
to -0.0718165785074234 -0.06219632923603366 -0.11162675172209053 -0.0351937639859581 0.043366577529907 0.009930439293384552 -0.06467465311288834 -0.49054917699310
26825 0.0883138177990527 0.025115570533466234 -0.43429619730738678 -0.2750214220903966 -0.0353405736386776 -0.0367852038227844 -0.095753252562561 0.1377663770673787
-0.2589706347405155 -0.2228063438036760 -0.1143641273954127 -0.33120733499527 -0.0591867561087036 -0.2348829507827758 0.2785216867937366 0.0513771644213195
17957484722137 0.24900034070014954 0.11283621791978455 -0.1255256626343536 -0.07017210125923157 0.00941756181418895 0.10308688732981682 0.07272831350564957 0.141226
with -0.0985988451363373 -0.2385247264729862 0.1175894141972064 0.06500596416130066 -0.0687106640954973 -0.2874845266342163 -0.0856295964479446 0.1696904450654
164905246347018895 0.21965739130973816 -0.3124369978904742 -0.102904126048084907 -0.0461577698588373 0.04010459035639945 0.2183129227161407 0.0417809662819091 0.22014
for -0.0885828405618668 0.0496838003300949 0.0377464219972778 0.06489270677478821 0.05439823531045 -0.112708864450065 -0.04672267660496819 0.1040973588824
1584 -0.21206684363112946 -0.039591915905475616 -0.6262766969906681 -0.2139615416527944 0.05573865562677383 -0.08018310368061066 0.13496974753836972 0.12017136662715202
covid-19 0.33330416673982324 0.01479049479856491 0.147060068098961833 -0.023025231392703056 0.191682571331787 -0.11672676154646683 0.08899730910759134 0.146422351817
55426 -0.04511909931898117 0.2939341962337494 0.24353934824466705 0.17898909747600555 -0.00947872167663956 -0.3241097927093506 -0.07146181166172028 -0.196622908115386
patients 0.04939527404561844 -0.1442393667353493 -0.1095365081378937 -0.25172389140905364 -0.08553808729410172 -0.0432542273368835 0.195867627859156 0.3668800294939
922166109085 -0.0716549913521688 0.0436086118221283 -0.02143552626831627 0.03729140389593545 0.03907901048660728 -0.0972571071789856 -0.0541148111224745 0.17745845
16
we 0.5449550747871999 0.12638379633426666 0.0373534739930373 0.08160431683063507 -0.11231604963541031 0.15202890336513552 -0.3664050400257106 0.1093474701046943
883780777454376 -0.6140636205673218 -0.09264766424894333 0.023603010922670364 -0.13203932344913483 0.09088529646396637 0.2482601147651672 0.026199063286185265 0.11704
2937
is -0.201849744273529 -0.0146817662048364891 -0.27701535820861 -0.3487420082922885 0.5759435891025103 0.1453063189983668 0.0859629511833109 0.42848418206596375
[UNK] 297199554443 0.04105306416749954 -0.1034683212679013 0.026820257306098

Word2Vec Visualization



Discussion

In these two graphs, we can find that in BERT result word mainly separate in two parts while Word2Vec result stay in one part. Also, in BERT result, there exists some nodes which stay very close. Although, Word2Vec result's nodes all exist in one part, the distance between two nodes is longer than BERT model. BERT result is more obvious, which may be one reason of playing a better role in word embedding.