**Chapter**     04. 자연어처리 (Natural Language Processing)

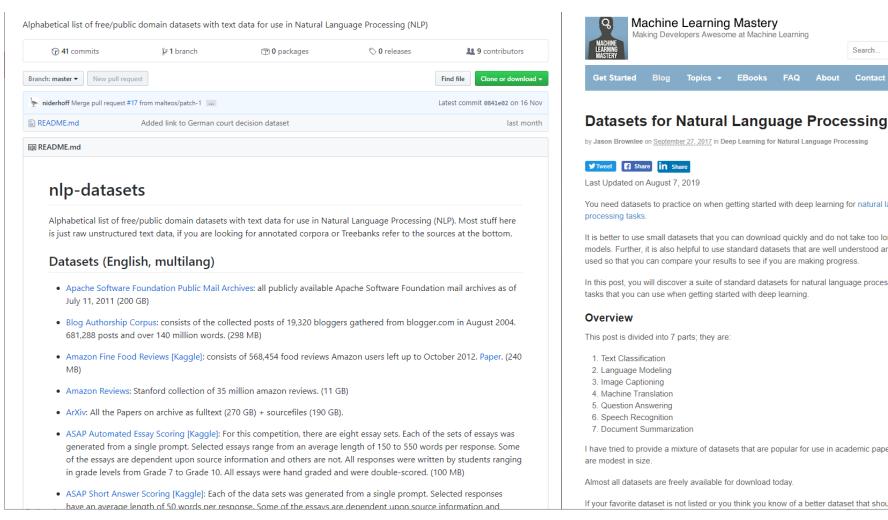# 응용하기 좋은 데이터셋 소개

# 데이터셋의 중요성

Awesome Deep Learning Algorithm

| Architecture | Weights | ← **Knowledge** | Learning Methods | Data |

아무리 잘 짜여진 딥러닝 알고리즘도, '지식'이 없으면 무용지물이다. 그 '지식'은 데이터셋으로 부터 나온다.

# 데이터셋 창고



https://github.com/niderhoff/nlp-datasets

https://machinelearningmastery.com/datasets-natural-language-processing/

# CoNLL Shared Task



https://www.conll.org/2019-shared-task

# Stanford Large Movie Review Dataset

## Large Movie Review Dataset

This is a dataset for binary sentiment classification containing substantially more data than previous benchmark datasets. We provide a set of 25,000 highly polar movie reviews for training, and 25,000 for testing. There is additional unlabeled data for use as well. Raw text and already processed bag of words formats are provided. See the README file contained in the release for more details.

Large Movie Review Dataset v1.0

When using this dataset, please cite our ACL 2011 paper [bib].

### Contact

For comments or questions on the dataset please contact Andrew Maas. As you publish papers using the dataset please notify us so we can post a link on this page.

### Publications Using the Dataset
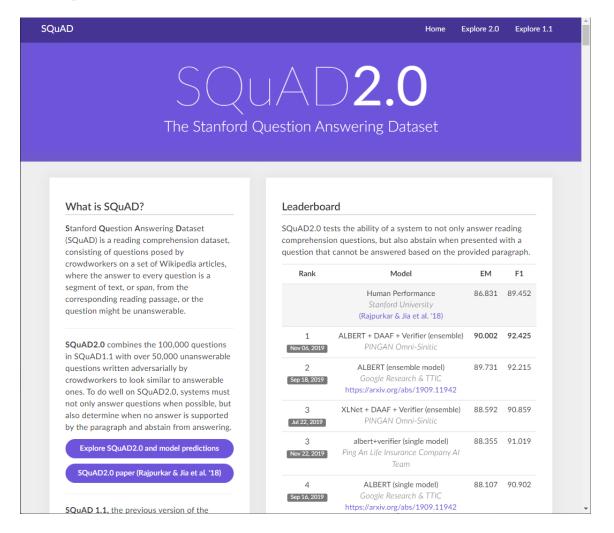
Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. (2011). Learning Word Vectors for Sentiment Analysis. *The 49th Annual Meeting of the Association for Computational Linguistics (ACL 2011).*
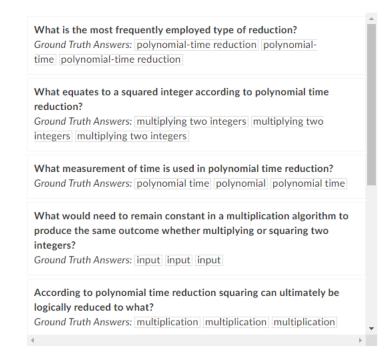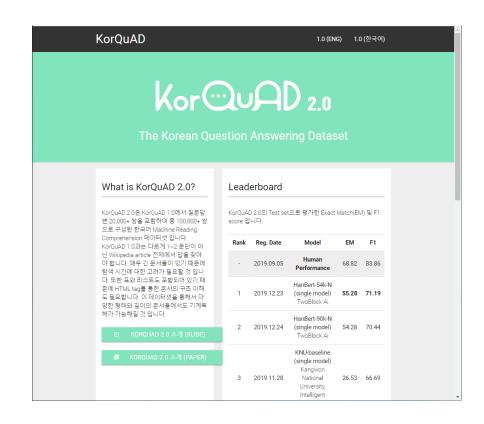
http://ai.stanford.edu/~amaas/data/sentiment/

# SQuAD2.0



https://rajpurkar.github.io/SQuAD-explorer/

# KorQuAD2.0 (2.1)



https://korquad.github.io/

# Naver Sentiment Movie Corpus

## Naver sentiment movie corpus v1.0

This is a movie review dataset in the Korean language. Reviews were scraped from Naver Movies.

The dataset construction is based on the method noted in Large movie review dataset from Maas et al., 2011.

### Data description

- Each file is consisted of three columns: `id`, `document`, `label`
  - `id` : The review id, provieded by Naver
  - `document` : The actual review
  - `label` : The sentiment class of the review. (0: negative, 1: positive)
  - Columns are delimited with tabs (i.e., `.tsv` format; but the file extension is `.txt` for easy access for novices)
- 200K reviews in total
  - `ratings.txt` : All 200K reviews
  - `ratings_test.txt` : 50K reviews held out for testing
  - `ratings_train.txt` : 150K reviews for training

### Characteristics

- All reviews are shorter than 140 characters
- Each sentiment class is sampled equally (i.e., random guess yields 50% accuracy)
  - 100K negative reviews (originally reviews of ratings 1-4)
  - 100K positive reviews (originally reviews of ratings 9-10)
  - Neutral reviews (originally reviews of ratings 5-8) are excluded

### Quick peek

```
$ head ratings_train.txt
id      document        label
9976970 아 더빙.. 진짜 짜증나네요 목소리          0
3819312 흠...포스터보고 초딩영화줄....오버연기조차 가볍지 않구나           1
10265843        너무재밓었다그래서보는것을추천한다        0
9045019 교도소 이야기구먼 ..솔직히 재미는 없다..평점 조정        0
6483659 사이몬페그의 익살스런 연기가 돋보였던 영화!스파이더맨에서 늙어보이기만 했던 커스틴 던스트가 너무나도 이뻐보였다  1
5403919 막 걸음마 뗀 3세부터 초등학교 1학년생인 8살용영화.ㅋㅋㅋ...별반개도 아까움.         0
7797314 원작의 긴장감을 제대로 살려내지못했다.     0
9443947 별 반개도 아깝다 욕나온다 이응경 길용우 연기생활이몇년인지..정말 발로해도 그것보단 낫겟다 납치.감금만반복반복..이들
7156791 액션이 없는데도 재미 있는 몇안되는 영화 1
```

https://github.com/e9t/nsmc