



BiliVista

Bilibili data analysis platform

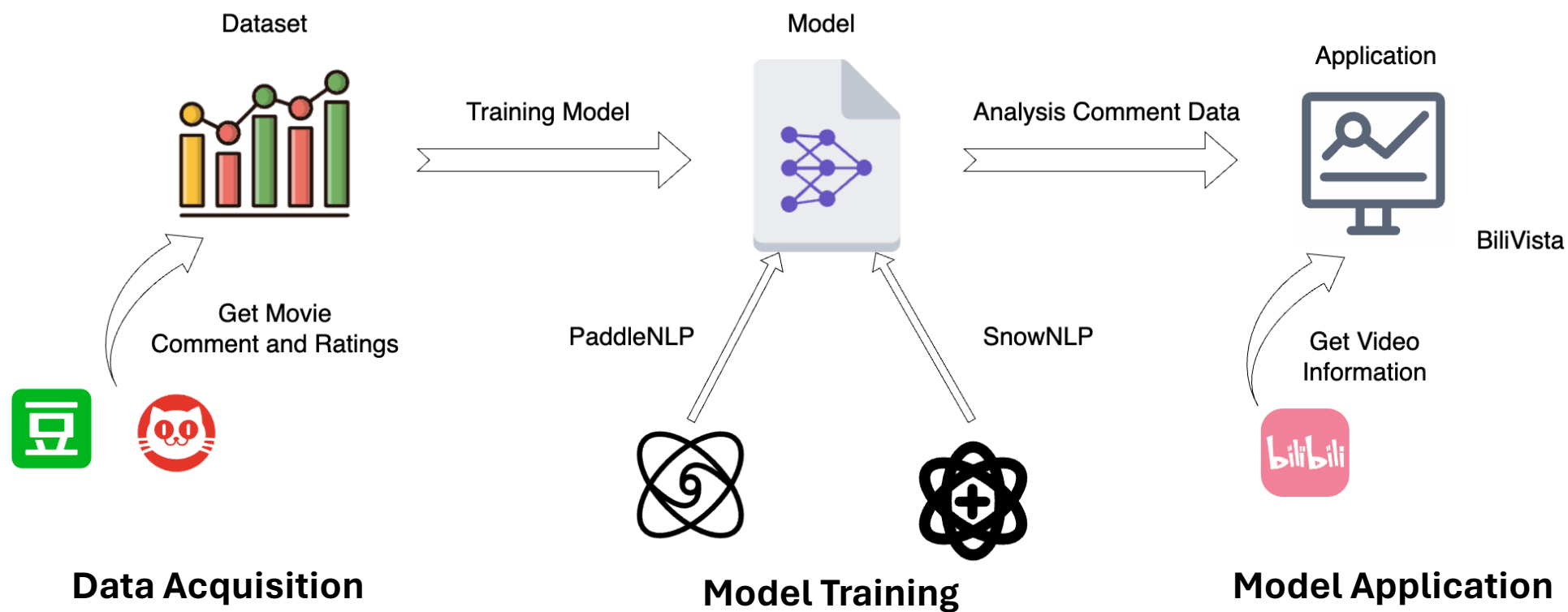
Emotion analysis based on deep learning



Project Overview

Data Acquisition

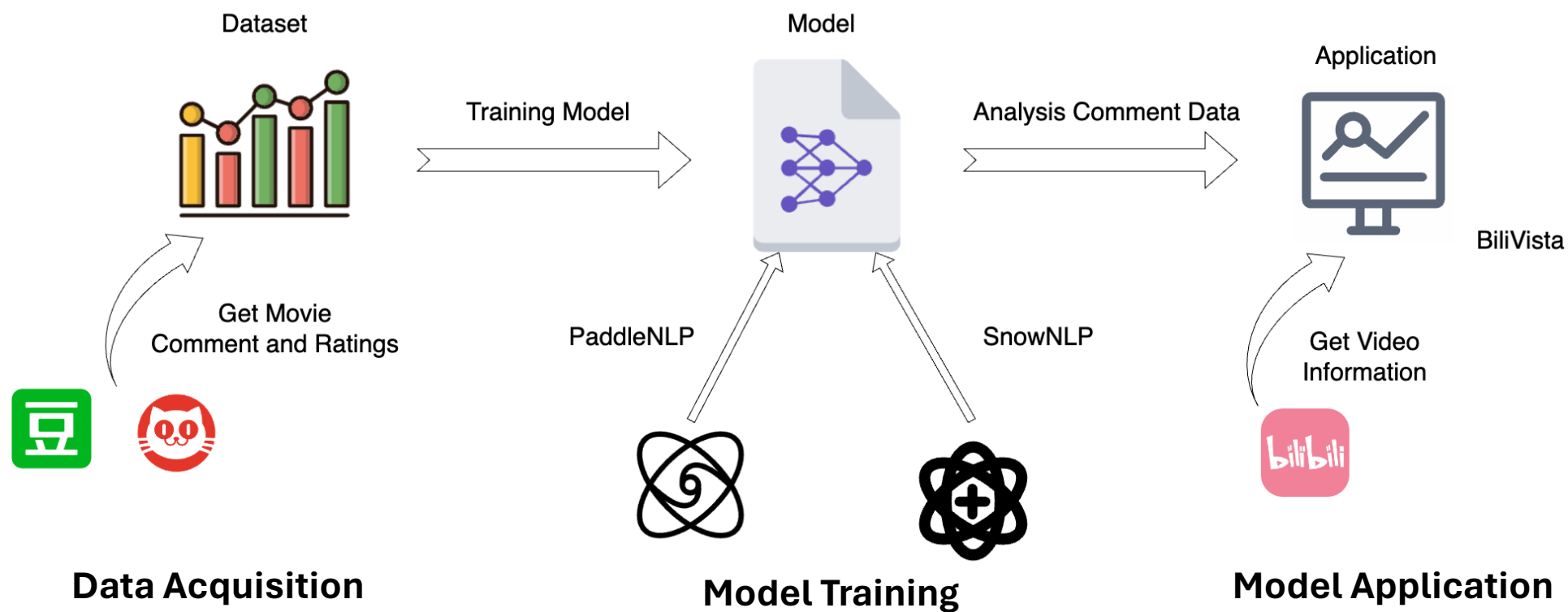
- **Crawling Movie Comments**
 - *Douban*: Collect comments from the `Douban` movie site.
 - *Maoyan*: Collect comments from the `Maoyan` movie site.
- **Data Preparation**
 - Perform data cleaning to ensure the quality and consistency of the collected comments.



Project Overview

Model Training

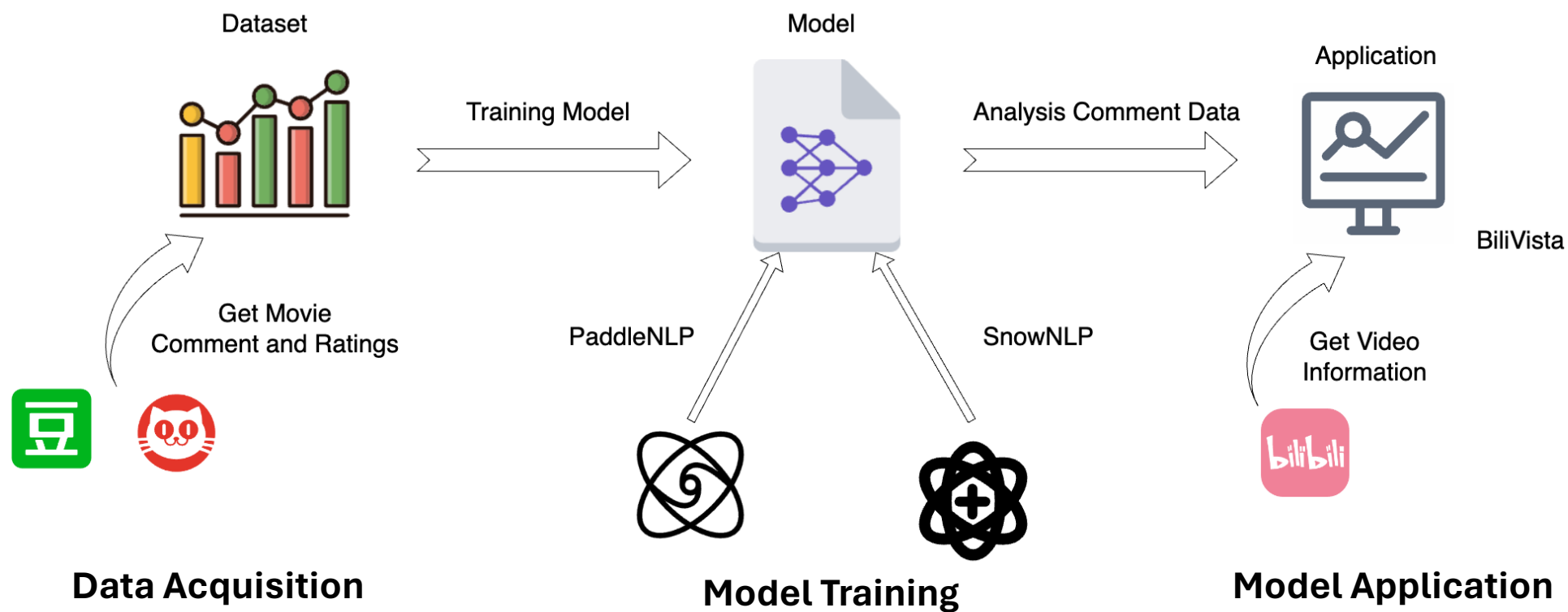
- **Machine Learning Model**
 - Utilize SnowNLP for sentiment analysis
- **Deep Learning Model**
 - Employ PaddleNLP to build and train deep learning models for advanced text analysis.



Project Overview

Model Application

- *Crawling Data from `Bilibili`*
 - *Implement real-time data crawling from `Bilibili`*
- *Backend Development*
 - *Develop the backend using the ****FastAPI**** Python*
- *Frontend Development*
 - *Implement a user-friendly interface to visualize and interact with the analysis results.*



Data Acquisition

Source of data

- *Douban*
- *Maoyan*

Why These Datasets?

- **Labeled Data:** This is **critical** as labeled data provides a **foundation** for training and evaluating our sentiment analysis models with **higher accuracy**.

Crawling Movie Comments

- *Douban:* Collect comments from the *Douban* movie site.
- *Maoyan:* Collect comments from the *Maoyan* movie site.

Leeloo 看过 ★★★★★ 2024-02-29 14:11:31 中国香港 5269 有用
沙丘美学算是被维伦纽瓦玩明白了。这次没想到的是被阿克南母星的建筑和审美追求折服了，无法用言语表达那种陌生而危险的美感：想象一下几万个伏地魔在人头攒动...查了一下原来作者就是有意将人们对阿克南家族的印象和苏联挂钩，所以我果然还是被苏式审美打动了orz

蕉叁叁 看过 ★★★★★ 2024-03-02 00:04:01 中国香港 3078 有用
太震撼了。这样毫无想象力的东西被人夸成史诗。hans zimmer的轰隆轰隆配乐如此刺耳。

简单.点点 Lv2
★★★★★ 3分 购票
烂片！（请相信这简短而又铿锵有力的评论）

风吹过的街道 Lv2
★★★★★ 10分 购票
努力活着 开心过好每一天

Data Acquisition

Data Preparation

- Perform data cleaning to ensure the quality and consistency of the collected comments.

Data Categorization

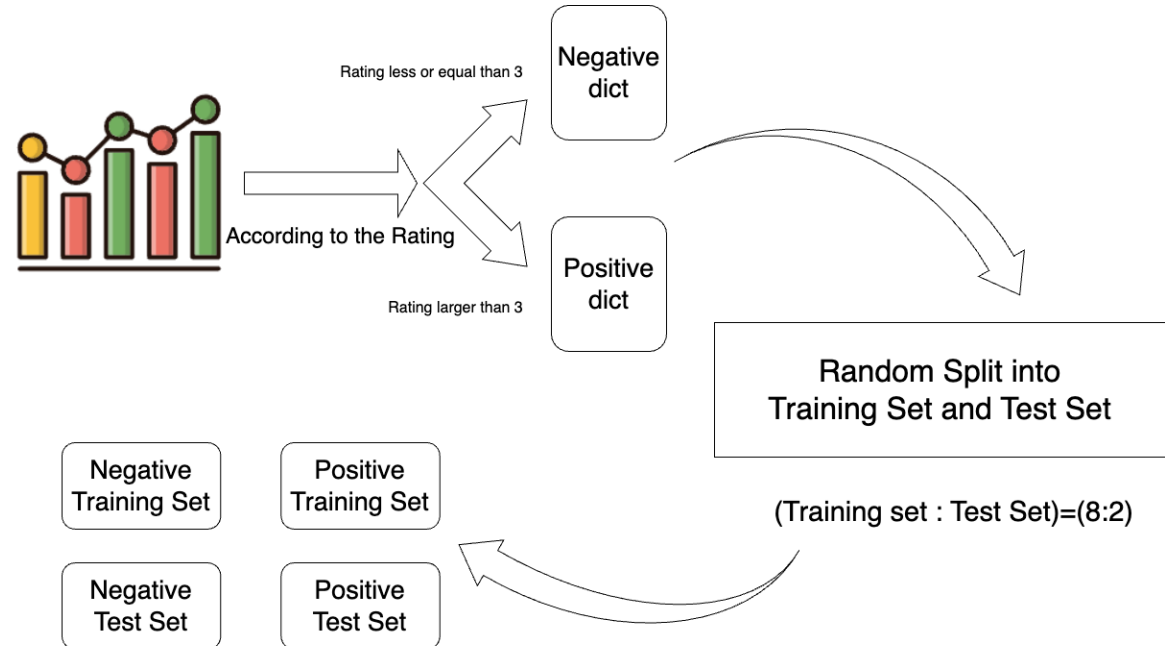
- **Negative dict:** Ratings less than or equal to 3
- **Positive dict:** Ratings greater than 3

Dataset Splitting

- Random Split

Training and Testing set

- Training set : Test set = **8:2**



Model Training

Machine Learning Model

- Utilize SnowNLP for sentiment analysis

SnowNLP

- SnowNLP is a library focused on natural language processing tasks for Chinese text, such as sentiment analysis and text processing.
- Uses the labeled training data to train the classifier using the [Naive Bayes algorithm](#).

Naive Bayes Algorithm

- Assumes that features are independent of each other
- Estimates probabilities based on the features and labels in the training dataset.

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

$$x_i = (x^1, \dots, x^n)$$

$$y_i = c_k \quad \text{and} \quad k = 1 \dots K$$

$$y = \operatorname{argmax}_{c_k} P(y = c_k) \prod_j P(x^j | y = c_k)$$

Model Training

Machine Learning Model

- Utilize SnowNLP for sentiment analysis

Naive Bayes Algorithm

- Assumes that features are independent of each other
- Estimates probabilities based on the features and labels in the training dataset.

Split Chinese characters and calculate the probability that the term appears in the set

Comments	Cut Sentence	Cut_Comments
很好, 揭露了现实打工人的心酸, 推荐	→	很好揭露了现实打工人的心酸推荐
还不错 有一些地方挺搞笑的!		还不错有一些地方挺搞笑的
非常搞笑, 值得推荐		非常搞笑, 值得推荐
太棒了, 太真实了, 推荐给领导们看看		太棒了太真实了推荐给领导们看看
搞笑喜剧, 挺现实的社会现象		搞笑喜剧挺现实的社会现象

Training Result

Method	Test Dataset Accuracy
SnowNLP	78.58%
KNN	78.32%

SnowNLP Model Evaluation

Practical application scenarios

- We actually using this model to annalysis the comment in other platform, it does not work well

Comment of Movie

- Crawl the Comment of the Movie 《Wandering Earth 2》 From the Bilibili

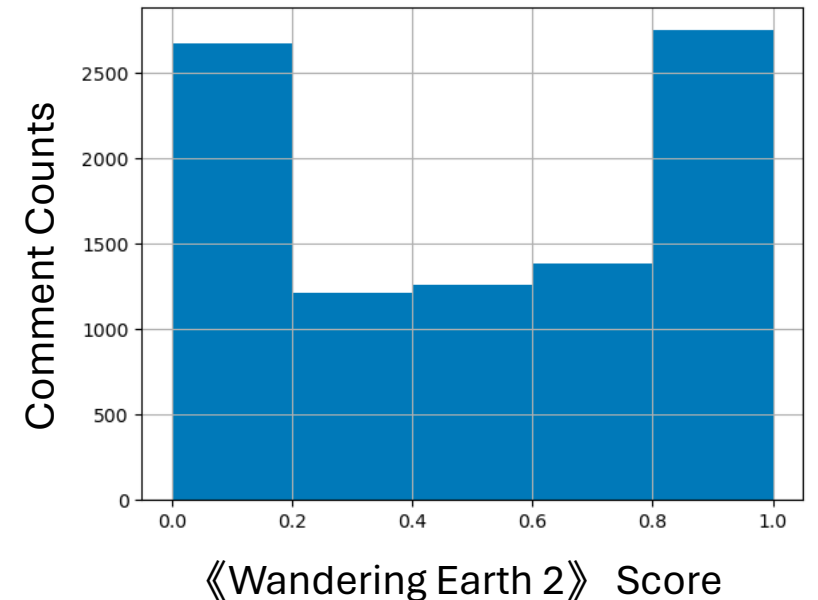
Estimation the comment score

- Get a comment sentiment score using the trained SnowNLP model.
- Average comment sentiment score is only 0.54.

(Assuming that the number of likes is the number of approvals)

Comment Score Distribution

- The distribution is polarized



SnowNLP Model Evaluation

Comment of Movie

- Crawl the Comment of the Movie 《Wandering Earth 2》 From the Bilibili

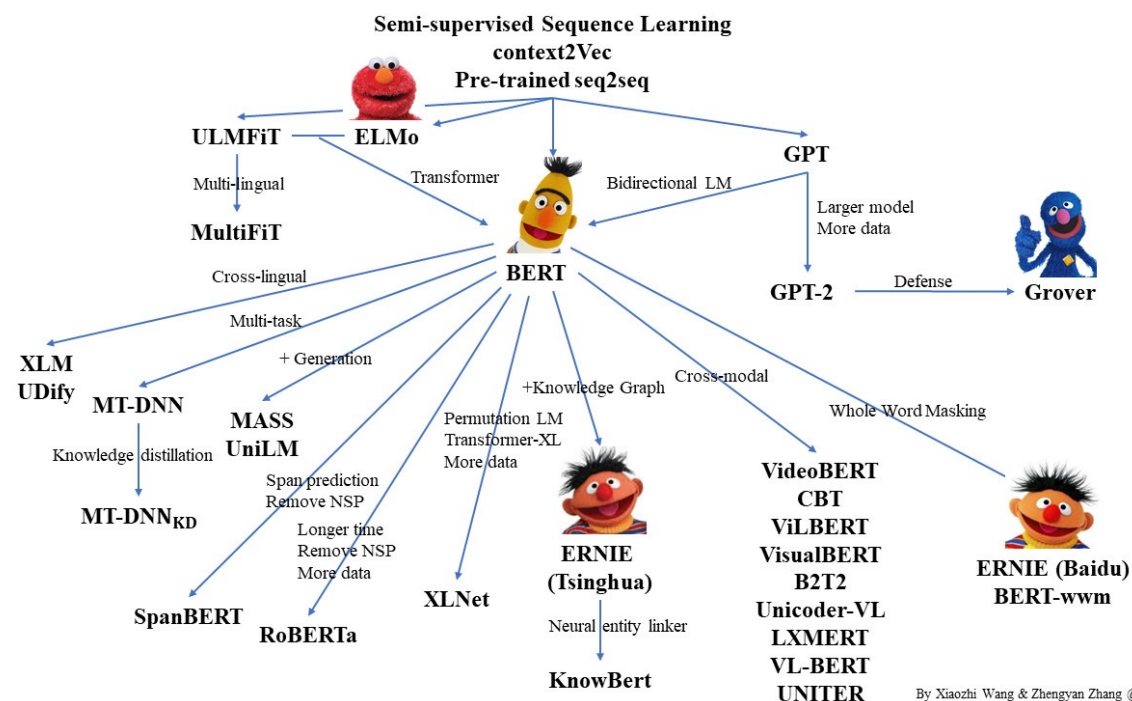
Comment	Like	Sentiment	Sentiment x Like
六公主给流浪地球2的颁奖词： 这是中国电影工业的一次全面跨越升级，以硬实力将中国科幻电影提升到前所未有的境界，7万多字原创剧本，2万多名工作人员，90多万平方米的置景总面积，历时1400余天的摄制，奉上一场2小时53的视觉盛宴。如此庞大精良的制作规模，造就了影片同名话题超11.8亿的网络关注度，收获40.23亿票房。 M大数据显示，影片传播指数达9.8，位列年度影片之首。 片中昂扬的中国精神，如同闪耀的星群，照耀着中国科幻电影前进的方向。	13621	4.14E-11	5.64E-07
其实导演郭帆不是半路出家，而是从小励志拍科幻片。 郭帆15岁时看了卡梅隆的《终结者2》，然后立志以后拍科幻片，他高考本来要考电影学院，但山东省没有招生导演系的，郭帆母亲也劝他考法律后当个政法委书记就行。 郭帆考上法律专业后想如果自己以后不奔着梦想去，等晚年躺病床、摇椅就特别后悔，所以他就觉得不管选什么专业，只要奔着哪个目标去，然后他大学也拍过电影短片，而且他学过法律学专业很适合工业化方面。 还有郭帆小时候画画很好，也拿过奖，有美术基础的。所以不要总半路出家、非科班、中途转行也能成功，搞得好像人家外行的行我也行，郭帆是自己本来就有这梦想、本来就有相关知识、而且他29岁时还考上北京电影学院管理系研究生。 中国科幻元年必定是1999年。	7597	0.000216	1.639879
郭帆在这一年高考，而且这一年的《科幻世界》压中了高考作文《假如记忆可以移植》，郭导看过并且受到启发拿了高分，同年也是大刘在科幻世界上开始首次投发文章《微观尽头》和《鲸歌》，次年就投发了《流浪地球》。同年高中生谢楠在1月的科幻世界发表了短篇奇想，页码43是吴京的生日4月3号。	6672	8.11E-10	5.41E-06

It turns out that the model does not work well in practice

Model Training – Deeping Learning

Pre-training model

- With the development of deep learning, the number of model parameters has increased rapidly, and in order to train these parameters, **larger data sets are needed to avoid overfitting**.
- Nowadays, studies have shown that **Pretrained Models (PTM)** based on large-scale unlabeled corpora can acquire generic language representations and **perform well** when **fine-tuned to downstream tasks**.
- **In addition**, pre-training models can avoid training models from scratch.



By Xiaozhi Wang & Zhengyan Zhang @THUNLP

Model Training

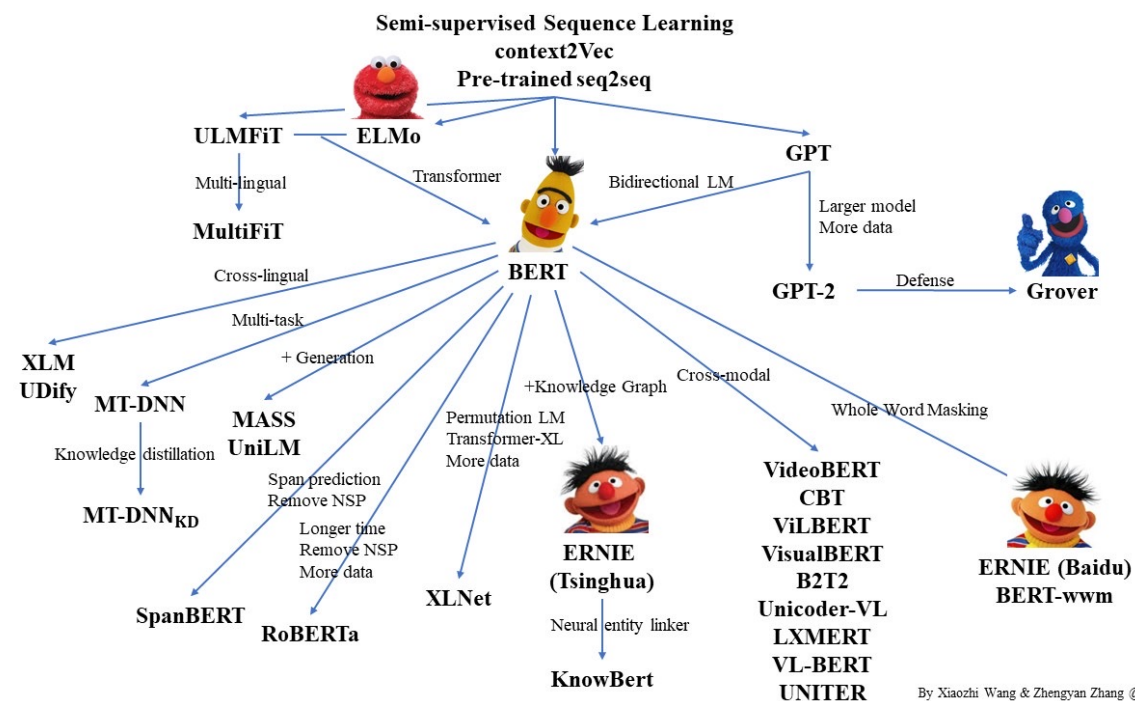
Pre-training model

- ERNIE (Like Bert-wwm)

BERT requires minimal architecture changes for a wide range of natural language processing applications.

Deep Learning Model

Employ PaddleNLP to build and train deep learning models for advanced text analysis.

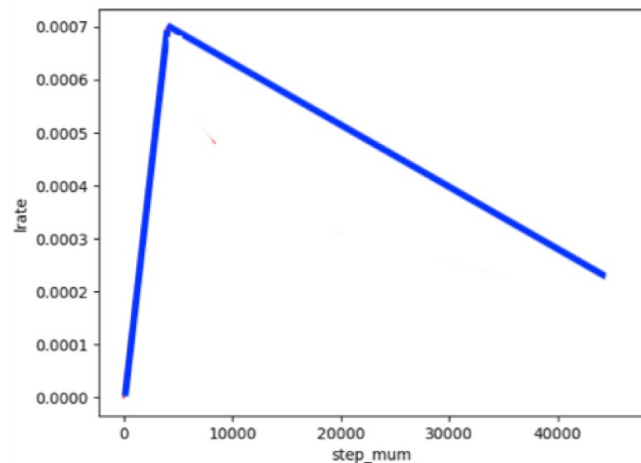


Model Training

Training Process

1. Remove a batch data from the dataloader
2. Feed batch data to the model for forward calculation
3. Pass forward calculation result to loss function to calculate loss. The forward calculation result is passed to the evaluation method, and the evaluation index is calculated.
4. Loss reverse return and update gradient. Repeat the above steps.

Each time an epoch is trained, the program will evaluate the effectiveness of the current model training.



Method	Test Dataset Accuracy
SnowNLP	78.58%
PaddleNLP	85.31%

PaddleNLP Model Evaluation

Practical application scenarios

- We actually using this model to annalysis the comment in other platform, it does not work well

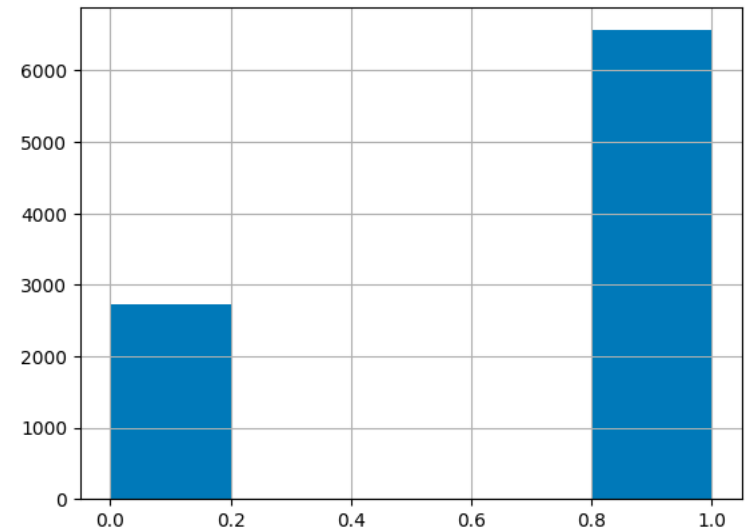
Comment of Movie

- Crawl the Comment of the Movie 《Wandering Earth 2》 From the Bilibili

Estimation the comment score

- Get a comment sentiment score using the trained SnowNLP model.
- Average comment sentiment score is 0.89.

(Assuming that the number of likes is the number of approvals)



《Wandering Earth 2》 Score

Model Application

Set the User-Agent and cookie information

- Set the User-Agent and cookie information
- Input the video bvid parameter to obtain the video.

Crawl video basic information

Through bvid access to the basic information of the video including:

- Title
- Author
- Reply count, Favorite count,
- Coin count, Share count
- ...

Crawling Data from *Bilibili*

- Implement real-time data crawling from *Bilibili*

```
1 HEADERS = {
2     'User-Agent': '',
3     'cookie':"",
4 }
5 url = 'https://api.bilibili.com/x/web-interface/view/detail'
6 params = {
7     'bvid': 'BV117411r7R1'
8 }
```

```
1 {
2     "BV1ct421u7gu": {
3         "title": "\u3010\u6708\u3011\u602a\u517d\u53f7 06",
4         "pic": "http://i1.hdslb.com/bfs/archive/95d2e0fcb4bcc9b5dfe334dff7b087595cce2e44.jpg",
5         "view": 1572285,
6         "danmaku": 10845,
7         "reply": 1030,
8         "favorite": 1320,
9         "coin": 5196,
10        "share": 352,
11        "like": 18884,
12        "author": "\u54d4\u54e9\u54d4\u54e9\u756a\u5267"
13    }
14 }
```

Model Application

Crawling Data from *Bilibili*

- Implement real-time data crawling from Bilibili

Crawl video comment content and danmu content

Through bvid access to the comment information of the video including:

- Comment
- Number of comments and likes

```
1 {
2   "comment": "\u6211\u6709\u4e00\u4e2a\u95ee",
3   "like": 3
4 }
```

Through bvid request screen XML file obtained through basic information:

- Time
- Timestamp
- Danmu text
- ...

```
1 {
2   "time": 4803.281,
3   "type": 1,
4   "font_size": 25,
5   "color": 16777215,
6   "timestamp": 1680401167,
7   "layer": 8,
8   "text": "\u4ed6\u7684\u59bb\u5b50\u4e5f\u662f\u6674\u5973"
9 },
```


Model Application

Backend Development

- Develop the backend using the **FastAPI Python**
- **FastAPI** is a modern, fast (high-performance) web framework for building APIs.

Why FastAPI?

Key Advantages:

- High Performance
- Rapid development
- Automatic interactive API documentation

Usage of 3 APIs:

- **video**: Fetch specific video information.
- **videos**: Gather statistics on all scraped videos.
- **newvideo**: Add **new** video data to the system.

Backend Architecture Overview

VideoInfo		^	
GET	/video/infos/	Get Video Info	v
GET	/video/attr/{attribute}/	Get Video Attribute	v
GET	/video/danmu/	Get Danmaku	v
GET	/video/comment/	Get Comment	v
GET	/video/isRunning/	Get Is Running	v
VideosInfo		^	
GET	/videos/infos/	Get Videos Infos	v
GET	/videos/counts/{type}/rank/	Get Rank By Type	v
GET	/videos/sentiment/rank/	Get Sentimentrank	v
GET	/videos/anyRunning/	Get Any Running	v
NewVideo		^	
GET	/newvideo/infos/	Get New Video Info	v
GET	/newvideo/comments/	Get New Video Comments	v
GET	/newvideo/danmu/	Get New Video Danmakus	v

Demo

Frontend Development

- *Implement a user-friendly interface to visualize and interact with the analysis results.*

Thanks