Research on Word Embedding Model for Automatic Generation of Novelty Retrieval Expression

Tingting Wang IMINZU University of China No. 27, Zhongguancun South Street, Haidian District, 100081 Beijing, China 1041987149@qq.com

Received June 2022; revised Sep 2022

Abstract: With the rapid development of Internet technology, text data is growing exponentially. How to effectively analyze and utilize these data and fully explore the value contained therein is the primary task of text big data analysis and processing. Text representation is an important work in the field of natural language processing, and how to better represent text semantics is an important cornerstone of practical applications in the field of natural language processing. The word embedding vector representation obtained by training can be considered to represent the word itself and its meaning. However, words have different meanings in different contexts, and polysemy is involved in different contexts. How to accurately express the word embedding vector to adapt to different contexts is also a hot topic of current research. At present, word embedding technology is still in its infancy, and there are still many problems worthy of further study. This paper makes a comparative study of the existing word embedding models, and selects the most suitable vector representation model for the automatic generation of novelty search formula to provide technical support for the scientific and technological novelty search system.

Key words: Natural language processing, Word-Embedding, Neural network, Novelty retrieval expression

1. Introduction. With the continuous development of modern society, data is becoming an important asset of individuals and enterprises, and the technology of analyzing and extracting the value of data has become the focus of attention in the era of big data. Users can use the technology of transforming natural language into structured processing to improve the efficiency and speed of mining effective information in structured data assets. At the same time, the application of this technology also lowers the threshold for non-technical users. This paper focuses on the research of Word-Embedding sub-task in the automatic generation technology of novelty search formula under the platform of novelty search system, and intends to select the most suitable model for novelty search formula. Word-Embedding model aims to map words into a vector space, each word has a dense vector of fixed length, similar words have similar vector representations, and polysemy can also be represented as multiple vectors in some models, which can facilitate quantitative text processing. Mining features between words and sentences in the text. The

representation of word vector has changed from the co-occurrence matrix and SVD decomposition based on statistics to the popular neural network based language model.

2. **Research on Word-Embedding Model.** Vector space model has been used in distributed semantics since the 1990s. At that time, many models for predicting word representations in continuous spaces had been developed, including Latent Semantic Analysis (LSA) and Latent Dirichlet Allocation (LDA). Bengio et al. [1] coined the term word embedding in 2003 and trained it jointly with model parameters in a natural language model. Collobert and Weston [2] showed the first practical application of pre-trained word-Embedding in 2008. Their landmark paper, A unified architecture for natural language processing, not only established word embedding as a useful tool for downstream tasks, but also introduced the neural network architecture that is now the basis of many methods.

But it was Mikolov [3-4] et al. Who created word2vec in 2013, a tool suite that allows seamless training and the use of pre-trained word-Embedding, that eventually made word-embedding popular. In 2014, Pennington [5] released a competitive set of pre-trained word-embedding GloVe, marking that word embedding have become mainstream. Word embedding is one of the successful applications of unsupervised learning. Their greatest benefit is undoubtedly that they do not require expensive manual annotation, but are derived from off-the-shelf big data sets that have never been annotated. Pre-trained word embedding can be used in downstream tasks that use only a small amount of annotated data. This paper will introduce the one-hot vector, word2vec, GloVe, ELMo and BERT vector, and analyze them combined with the neural network language model to get an efficient and fast automatic generation model of retrieval.

3. Word-Embedding model

One-hot Encoding 3.1. One-hot encoding is also known as one-bit efficient encoding. It is processed by occupying the N-bit status register, but only one of the N bits is valid, that is, 1, and the rest are all 0. This requires mapping the classification values to binary integer values, with each integer value labeled as a binary vector. Detailed explanation of the coding process: first determine the set of coding objects, and then determine the elements in the set. For example, the phrase "hello world" has 27 categories (space+26 lowercase letters). In this example, there are 11 elements, each of which has 27 features, which are converted into a binary vector representation as shown in Figure 1. (1) The 27 features are first integer encoded: a--0, B--1, C--2.., Z--25, space--26, (2) The 27 features are arranged from front to back according to the size of the integer code.

Categorical variables

-																													
7	а	b	с	d	е	f	g	h	i	i	k	1	m	n	0	p	a	r	s	t	u	v	w	x	v	z	容		hinan worto
h	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	>	Dinary vecto
е	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0		
空	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1		
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0		
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0		
r	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0		
1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
d	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		

FIGURE 1. One-hot Encoding

One-hot encoding is intuitive, with each element corresponding to only one feature, making it easy to calculate the loss function or accuracy. It also has two serious disadvantages: (1) The length of the dictionary is the length of the representation of each dimensional vector, and only one position in each dimensional vector is 1, which wastes memory and is not conducive to calculation. (2) One-hot coding is a bag-of-words model, whose vector matrix is equivalent to numbering each word, and the resulting features are sparse and discrete. However, the words in the text are related to each other, which does not take into account the semantic problems of the co One-hot converts words into discrete individual symbols, and the hidden layer is a linear unit without activation functions. Using Soft-max Function for normalization, the dimensions of the output layer are the same as those of the input layer. Word2vec was proposed after the emergence of one-hot. The Purpose of Word2vec is to transform natural language into Dense Vector, and to transform sentences and words in natural language into Dense Vector through the model of Continues Bag of Words. The two models of CBOW and Skip-gram are transformed into dense vectors that computers can understand and calculate, and the dense vectors we want to obtain are actually the output units of the hidden layer. Word2vec is an intermediate result (weight matrix) generated by neural network learning, and its downstream task is to serve as the corpus input of neural probabilistic language model. Hierarchical Soft-max or Negative Sampling is used in the neural network learning process to reduce complexity. The CBOW/skip-gram model is shown in Figure 2.





In the figure, the former word represents the word, and the latter word represents the word. The running window of the CBOW model in the figure 2, which is predicted by the upper and lower words of the target word. The prediction process is as follows: first, the dimension of the word vector is set, here for example, dimension d, and then all words are randomly initialized as the word vector with the set dimension d. Then the vector of the target word is predicted by the hidden layer vector obtained by encoding the word vectors of the two words, and above and below the target word. Finally, it can be seen from the figure that in the bag-of-words model, he simply adds the vectors of the predicted target words, and then does Soft-max normalization to get the target words. For example, if there are V words in a vocabulary set, the D-dimensional hidden layer vector is multiplied by the W matrix and converted into a V-dimensional vector for soft-max classification. Generally, the number of vocabulary is very large, and the number of occurrences of different words varies greatly. In the experiment, soft-max and negative sampling are generally used to optimize the problem of large amount of calculation.

The Skip-gram model predicts the number of words in the window size of the context through the target word, that is, the target word is input and mapped to the hidden layer vector, and the upper and lower words of the target word are predicted according to the vector of the target words. The problem of unbalanced vocabularies and samples is handled in the same way as CBOW.

GloVe model 3.2. GloVe model is a Word2vec model that integrates the overall information, and it can also be considered as an improved Word2vec model. From the previous chapter, it can be seen that CBOW and Skip-gram models only consider the local information of words in the window, but lose the semantic information of words outside the window, so GloVe uses co-occurrence matrix to consider the fusion of local information and global information. The title of this thesis is used to introduce the running window of GloVe. The size of the window is set to 5, that is, the two words before and after the head word.

Window	Center	Window content					
label	words						
0	Research on	Research on, Word-Embedding, Model					
1	Word-	Research on, Word-Embedding, Model, for					
	Embedding						
2	Model	Research on, Word-Embedding, Model, for,					
		Automatic					
3	for	Word-Embedding, Model, for, Automatic,					
		Generation					
4	Automatic	Model, for, Automatic, Generation, of					
5	Generation	for, Automatic, Generation, of, Novelty					
6	of	Automatic, Generation, of, Novelty, Retrieval					
7	Novelty	Generation, of, Novelty, Retrieval, Expression					
8	Retrieval	of, Novelty, Retrieval, Expression					
9	Expression	Novelty, Retrieval, Expression					

TABLE 1. Window content

In order to better express the contents in the table with vectors, the following symbols are used to analyze the rules:

 $X_{i,j}$: Obtained by traversing all the window content statistics, $X_{i,j}$ representing the number of times J appears in the window content with I as the central word.

 X_i : Count how many other words the word I is the head of. $X_i = \sum_{j=1}^N X_{i,j}$.

 $P_{i,k}$: The proportion of word K near word I to all nearby words of word I. $P_{i,k} = \frac{X_{i,k}}{X_i}$. Ratio_{*i*,*i*,*k*}: proportion of the proportion of word K occurring near word I and word $J,Ratio_{i,j,k} = \frac{P_{i,k}}{P_{j,k}}$.

The following table is obtained by statistical analysis of the values of Ratio_{i ik}:

TABLE 2. $Rutio_{i,j,k}$ Conclation scores with words							
Ratio	Words I K related	Words J, K not					
i,j,k	words J, K related	related					
Word I, K correlation	approaches 1	Very large					
Word I, K correlation	Very small	approaches 1					

TABLE 2. $Ratio_{i,i,k}$ Correlation scores with words

ELMo Model 2.4. Although GloVe incorporates local information and global information, it has not yet solved an existing problem, that is, polysemy. Words have different meanings in different sentences. For example, "bank" can refer to either a bank or a riverbank. In Chinese, it is even more so. For example, "apple" can refer to either a brand in electronic products, a mobile phone, a tablet or a computer, or an apple in the fruit world. In the GloVe model, these two words represent the same vector in different contexts. ELMo [6] solved this problem by extracting semantic and syntactic features from LSTM.

$$P(t_{1}, t_{2}, \dots, t_{N}) = \prod_{k=1}^{N} p(t_{k}/t_{1}, t_{2}, \dots, t_{k-1})$$
(1)

$$P(t_{1}, t_{2}, \dots, t_{N} = \prod_{k=1}^{N} p(t_{k}/t_{k+1}, t_{k+2}, \dots, t_{N})$$
(2)

In the forward process of Formula 1, the word from the first word to the (K-1)th word is used to predict the kth word, and in the backward process of Formula 2, the word from the (K+1)th word to the Nth word is used to predict the kth word in the reverse direction. ELMo predicts the target word through a double-layer bidirectional recurrent neural network RNN, as shown in fig. 3, where E represents the input of word features. The language model is trained by extracting syntactic features and semantic features using two-layer LSTM, and then the regression processing is performed using softmax.

Specifically, words are predicted in both directions, and the specific encoding method can be expressed by LSTM as:

$$\sum_{k=1}^{N} (logp(t_{k}/t_{1}, t_{2}, \cdots, t_{k-1}; \Theta_{x}, \Theta_{LSTM}, \Theta_{s}) + (logp(t_{k}/t_{t+1}, t_{k+2}, \cdots, t_{k}; \Theta_{x}, \Theta_{LSTM}, \Theta_{s})))$$
(3)

 Θ_x is the vector representation of Token and Θ_s is the parameter of softmax.



FIGURE 3. ELMo model

Bidirectional LSTM concatenates multiple layers, and the vector output of each hidden layer is used as the vector input of the next layer. For the word E_1 , the two-layer bidirectional LSTM language model as shown in Figure 3 can obtain two vector representations at each layer, and then +1 is the vector representation of the Token output layer, and the five vectors are combined to obtain the vector representation of each Token of ELMo. Then for the kth word, each Token will have 2L+1 vector representations.

$$\mathsf{ELMo}_{k}^{task} = E(R_{\chi}; \Theta^{task}) = r^{task} \sum_{j=0}^{L} \mathsf{s}_{j}^{task} h_{k,j}^{LM}$$
(4)

 $h_{k,j}^{LM}$ is the hidden layer output for the jth layer of Token K, and represents the vector representation of the input layer Token when J = 0, s_j^{task} is a normalization of soft-max and r^{task} is a scaling parameter, representing the proportion of each feature, which allows the target model to scale the vector of ELMo. ELMo solves the problem of polysemy, and the part of speech is also corresponding, but it still has two problems: ELMo's ability of using LSTM as language model to extract features is weaker than Transformer proposed by Google; LSTM, even BiLSTM, only considers unilateral data and sums them at the loss function, which is weak in the ability of bidirectional fusion of features. As a sequential model, the parallel computing ability is poor.

BERT model 2.5. BERT [7] has a wide range of applications, and has designed extremely simple downstream interfaces for sequence labeling, classification tasks, and sentence relation judgment tasks in the four major tasks of natural language processing. Figure 4 below is the interface designed for the downstream task in the BERT paper. BERT can be regarded as a black box, which works in a similar way to ELMo. It is a Pre-trained in the current large-scale corpus, and then the downstream task is input to compare the lightweight Fine-tuning. Downstream tasks can be put into BERT's Fine-tuning model with some slight modifications. It can be seen that for sentence-relational tasks, as long as the start and end symbols of sentences are added,

and the division symbols are added between sentences. Then the output of the last position of the BERT model is connected to a softmax classifier. For the model of sequence labeling, it is also enough to add the start and end symbols, and add a linear classifier to the output of each position of the final BERT.



FIGURE 4. BERT model

The subsequent BERT uses Transformer for encoding. The whole Transformer structure is divided into two parts as shown in Figure 5 below. One side is Encoder, and the other side is Decoder. The reason why Transformer will replace recurrent neural network and convolution neural network as the mainstream coding method of natural language processing is that convolution neural network extracts local features, but for text data, it ignores long-distance dependence, and the coding ability of CNN in text is weaker than that of RNN, while RNN is a sequence model with poor parallelism. It is slow to compute and can only consider information in one direction. Transformer can comprehensively consider the information in two directions, and comprehensively consider the context features in both directions when predicting words.



FIGURE 5. Transformer model

By looking at the block ENCODER # 1, for example, the input is a vector of a line of sentences with its position vector encoded and then passed through a Self-Attention. Self-Attention first maps the embedding of each Token to three vectors, Query (Q), Key (K), and Value (V), through three different matrices. Figure 6 below assumes that the length of the sentence is only 2 tokens. Since the three W conversion matrices of Q, K, and V share parameters at different positions, matrix operations can be used. The Q and K vectors are used to calculate the weight parameters of Attention. The calculated weight parameters are multiplied by the V vector, and then weighted and summed to obtain the vector of each token. The calculation process is as follows: For example, to calculate the vector of position 1, first, Q_1 is multiplied by K_1 to K_n points to obtain n scalars. This scalar is divided by $\sqrt{d^k}$ and d^k is the dimension of the vector of K. After softmax normalization processing, the normalized weight is correspondingly multiplied by V_1 to V_n , and the Token vector of position 1 is obtained by summation.

Self-Attention will also superimpose H independent Self-Attention mechanisms in parallel. The Token at each position will have H vectors. The H vectors will be spliced and a linear transformation will be made to get the final vector. H is called the head number of Attention. Because this Self-Attention structure will ignore the location information, unlike CNN and RNN, which will naturally encode the location information. The location information is very important for the sequence text. Add the location code at the output, and then make a standardization, corresponding to Add & Normalize in Figure 5. The position code here is calculated as follows:

$$\mathsf{PE}_{(pos\,2i)} = \sin(pos/10000^{2i/d_{model}}) \tag{5}$$

$$\mathsf{PE}_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{model}}) \tag{6}$$

The position of each POS is represented by a d-dimensional vector, the even-numbered position of the vector is calculated by a sine sin function, the odd-numbered position is calculated by a cosine cos function, and a value between -1 and 1 is obtained; the trigonometric function is used because the characteristic that the

sum and difference transformation of the trigonometric function can be linearly transformed is utilized, After a Feed Forward neural network and Add & Normalize, an encoding is completed and passed to the next encoding block. The structure and operation mechanism of Decoder are basically the same as those of Encoder. The difference is that the K and V vectors at the top level of Encoder will be passed to the Encoder-Decoder Attention component of each Block of Decoder. Finally, the final result is obtained through a linear transformation and the softmax classifier.

3. **Neural network language model.** The essence of deep learning lies in vector representation, but natural language processing has not done the most fundamental representation learning---Word-Embedding and language model well, and no model can achieve excellent results alone. Combining the existing models as functional modules to construct a more complex model is the mainstream idea to solve the problem of automatic generation of novelty retrieval.

From the beginning of BP neural network, neural network model began to get a lot of attention. Convolution Neural Network (CNN) is an efficient feature extraction method developed in recent years and has made significant breakthroughs in many tasks. CNN is a convolution neural network, which is used for image processing and extracting high-latitude features. Then Recurrent Neural Network [9] (RNN) appears. RNN has the problem of dimension explosion and gradient disappearance. Therefore, the Long Short-Term Memory (LSTM) model proposed by Hochreiter and Schmidhuber et al. [10] is improved. In 2014, K. Cho [11] proposed a gated recurrent unit network (GRU), which is another RNN gating architecture that has attracted attention after LSTM. GRU is a variant of LSTM with fast speed. In 2013-2015, Y. Bahdanau, etc. [12] proposed a series of RNN algorithms such as encoder-decoder and self-attention layer, and applied them to machine translation problems [13]. Seq2seq is used to solve the problem of unfixed length of input and output, and attention mechanism is used to solve the problem of information overload.

The Seq2seq model is also known as the Encoder-Decoder model. In deep learning, the encoder-decoder framework is a classical scheme to deal with the problem of sequence-to-sequence mapping. Another output sequence y is generated from one input sequence X. Seq2seq model can be applied in many aspects, such as machine translation, document extraction, question answering system and so on. The encoder is to encode the input sequence into a vector of fixed length, and the decoder is to convert this fixed vector into a sequence output again. The Encoder-Decoder model diagram is shown in Figure 6, where the decoder and encoder are not fixed, and the optional models are CNN/RNN/Bi-RNN/GRU/LSTM, etc. Encoder and Decoder can be freely combined from these models.



FIGURE 6. Encoder-Decoder architecture diagram

The loss function during the training of the Encoder-Decoder model generally uses the cross-entropy loss, that is, to maximize the probability of generating the annotated text in the training data, and $P(y_{i,k})|y_{(1,2,\dots,i-1)})$ represents the probability of Decoder generating a word with an ID of K in the ith step, so given an annotated text, the logarithmic probability of being generated by the model is:

$$P(y_1, ..., y_m) = -\log \prod_{i=1}^{m} p(y_{i,k}) = \sum_{i=1}^{m} \log (y_{i,k})$$
(7)

Encoder and Decoder commonly use RNN-like units such as Vanilla RNN, GRU, and LSTM. The RNN receives the current input and the hidden state from the previous input at each step through the self-loop, that is, the hidden state at the current time is determined by the state at the previous time and the current time input, which is expressed by the formula:

$$h_t = f(h_{t-1}, h_t) \tag{8}$$

The information of each hidden layer is combined to generate a semantic vector of

$$C = q(h_1, h_2, h_3, \dots h_{T_X})$$
(9)

A simple way is to take the last hidden layer as the semantic vector C, which is

$$C = q(h_1, h_2, h_3, \dots h_{T_X}) = h_{T_X}$$
(10)

Decoding is done by combining a given semantic vector C with the output sequence $y_1, y_2, y_3, \dots, y_t$ to predict the output of the next word y_{t-1} , which is

$$y_t = argmaxP(y_t) = \prod_{t=1}^{T} P(y_t | \{y_1, ..., y_{t-1}\}, C)$$
(11)

In the formula, C represents the semantic vector, and yt-1 is the output representation of the previous time period, which will be used as the input at this time again.

Although the Encoder-Decoder model is classical, its limitation is that encoding and decoding are only related by a fixed-length semantic vector C, which has two defects. One is that the fixed-length semantic vector can not fully express the input information, and the other is that the content input first will be covered and diluted by the content input later.



FIGURE 7. Seq2seq attention architecture diagram

In the process of generating search formula for novelty search, it is necessary to summarize and analyze the generation strategy of search formula for novelty search, the text vector representation model and the text generation model. This chapter introduces the word embedding model and the generation model used in the automatic generation model of search formula, which is intended to provide technical support for the realization of automatic generation function of search formula, and provide reference for the improvement of some models in this paper.

4. Characteristic acquisition of scientific and technological literature

Literature acquisition and preprocessing 4.1. Due to the scarcity of the orders provided to the novelty retrieval system platform at present, it can not provide sufficient data for this experiment, so the abstract of scientific literature is selected as the novelty retrieval point. In this paper, the scientific literature in the fields of computer, information science, chemistry and pharmacy are selected as the experimental objects, and the experimental data are retrieved from Wan-Fang Database Knowledge Service Platform. The retrieval strategy is to limit the abstract structure to "purpose: method: result and conclusion", and the time is from 2015 to 2020. The specific number of relevant documents retrieved is shown in Figure 8.



FIGURE 8. Distribution of Literature Quantity

It can be seen from the figure that the number of literature in the field of pharmacy

is the largest, and the number of literature in the field of information science is the smallest. In order to obtain the literature source of the specification, remove the irrelevant literature. It is necessary to preprocess the acquired scientific and technological literature. The specific process is shown in the figure below. Firstly, data collection and preprocessing are carried out, including merging duplicate items, deleting missing items, and extracting "title" and "abstract" fields to form a corpus to be analyzed and processed. Then the corpus is processed by word segmentation, stop word filtering and other steps to form the data set needed for the final experiment.

Text word segmentation is to generate a sequence feature document by processing the text. Chinese word segmentation is a unique concept in the process of Chinese processing. English text has obvious spaces to distinguish the boundaries between words, while the boundaries between words in Chinese text are not just simple punctuation marks. For this reason, the Chinese word segmentation is an important step in the preprocessing, but the Chinese word segmentation technology is not complicated. The following figure shows the word segmentation results of processing the text of the novelty search point using the appropriate word segmentation model.



FIGURE 9. Example of Text Segmentation

Stop word filtering is also a key step in text processing, deleting words such as "joint", "determination", "based on" that have no clear meaning but will always appear. These words are essential in the process of text reading and comprehension, but in the process of building the model, it has been assumed that words are interchangeable, so these words are of no value to the model of this paper. During the experiment, the text is processed by using the stop word list defined by ourselves. The following figure is an example of matching results using regular expressions after removing stop words.

Tc99 GSA 肝脏 储备 功能 立体 定量 评估 精准 肝 切除 三维重建 技术 联合 持久美 蓝 染色法 精准 肝 切除术 精准 肝 切除 理念 核心 最大化 病灶 去除 最大化 脏器 最小化 创伤 侵袭 理念 肝脏 外科 实践 精准 肝脏 外科 IGD13CNMR 技术 天然 产物 活性 成分 活性 成分 组 分子结构 解析 优势 IGD13CNMR 偶联 指纹 图谱 技术 植物 源 中药 植物 源中 兽药 植物 源 农药 基础 质量 控制 60 例 老年 患者 择期 行 全身 麻醉 脑肿瘤 手术 随机 分为 右美托咪啶 组 对照组 监测 记录 麻醉 诱导 时间段 平均 动脉 压 心率 变化 结果表明 Dex 抑制 气管 插管 MAP HR 升高 血流 动力学 稳定 术后 恢复 作用 腹腔镜 辅助 腋下 切口 Ivor Lewis 经胸 颈部 机械 吻合术 腹腔镜 经腹 食管 裂孔 清扫 纵隔 淋巴结 超声 刀经 胸 径 路 清扫 颈 胸 交界 颈 段 食管 旁 淋巴 组织 相结合 三野 淋巴结 清扫 手术 治疗 中段 食管癌 蛹 虫草 玛卡 黄精 枸杞 原料 辅以 菊粉 调制 成 增强体质 耐受 疲劳 一款 饮料 顶空 气 相色谱 质谱法 测定 蜂蜜 57 挥发性 有机溶剂 烷烃 类 芳香烃 类 醇类 酮类 酯类 醚 类 残留量 分析方法 添加 silicalite 晶种 控制 聚乙二醇 浓度 ZSM 粒径 可控 合成 ZSM 晶粒 尺寸 60 nm 14 um 之间 均匀 变 14. 顺丁橡胶 基体 材料 天然橡胶 丁苯橡胶 改性剂 Si69 Si75 偶联剂 制备 高耐磨 低 硬度 抗湿 滑性 新型 橡胶 外底 材料 TiO2 载体 稀土金属 Ce La Cu Mo Zr Fe Co Ni 三种 三种 过渡 金属 改性 掺杂 替换 传统 有毒 活性 组分 V205 构建 绿色 环保型 钒 纳米 复合 氧化物 CeLaCuMoNiOx TiO2 CeLaCuZrCoMoFe0x Ti02 CeLaCuFeNiOx TiO2 后以 四氯化碳 正 硅酸 乙酯 前 躯体 化合物 化学 液相 沉积 CLD 方法 复合 氧化物 定向 碳 硅 沉积 催化裂化 FCC 再生 烟气 洗择性 催化 还原 SCR 脱硝 催化剂 以一 氯乙酸 原料 DMF POCl3 HPF6 合成 氯 双 二甲基 氨基 三亚 甲六 氟 磷酸盐 含铅 砷 元素 乳状物 标准 物 制备 利培酮 口服 溶液 处方 利培酮 200g 酒石酸 1.0 氢氧化钠 100g 苯甲酸 200g 纯化 水加 200L 自动进样 器 通道 一体化 设计 蠕动泵 进样 多通道 16 通道 流动 注射 分析仪 牙周 干预 COPD 患者 辅助 治疗 辅助 效果 COPD 患者 肺 功能 指标 降低 COPD 患者 急性 发作 频率 COPD 患者

FIGURE 10. an example of a text removal stop word.

Feature extraction experiment 4.2.

One-hot encoding 4.2.1. One-hot coding is simple and easy to use, and is a commonly used method in the text feature extraction process. The flow of the extraction process is as follows: Firstly, each sentence in the corpus is divided into words and numbered, and the words in each sentence are matched with the numbered words, if they are matched, they are 1, and if they are not matched, they are 0.

Word2vec feature extraction 4.2.2. The process of Word2vec feature extraction is shown in the figure below:



FIGURE 11. Word2vec feature extraction flow

(1) Obtaining the preprocessed literature, searching data according to the pharmaceutical field selected in the experiment, and constructing a text corpus;

(2) The Word2vec model of the training corpus, in the training process, the parameters are set as vector_size, that is, the dimension of the feature vector is 300, min_count, the minimum word frequency of the words participating in the training is 1, windows_size means that the size of the training window is set as 15, and DM training algorithm is set as 1. The number of ITER iterations is set to 500.

(3) Vectorized representation of text, that is, a word is represented as a vector.

(4) Calculate the cosine values of the two vectors. The range of the cosine values is [-1, 1]. A value approaching 1 means that the two vectors are close in direction; a value approaching -1 means that the two vectors are opposite in direction; and a value approaching 0 means that they are nearly orthogonal.

GloVe feature extraction 4.2.3. The purpose of GloVe is to make the words vector contain the whole semantic information as much as possible. The process of feature extraction is firstly to construct the co-occurrence matrix based on the dictionary, and then to operate the vector based on the co-occurrence matrix by using Euclidean

distance or cosine similarity, and to construct the loss function. The input is a corpus and the output is a word vector. The process is shown in the following figure:



FIGURE 12. Flow chart of GloVe feature extraction

In the experiment, we need to use Linux or OS system. First, we need to download the GloVe code. The "src" folder contains four main running files. By running these four files, we can generate the vectors we need. Firstly, the vocab_count. C is executed, and a dictionary, i.e., a vocab. Txt, is generated by traversing the corpus; the "cooccur.c" file is executed, and a co-occurrence matrix is established from the corpus. It can be seen from fig.14 that the training of GloVe is performed on the basis of the co-occurrence matrix; Then, a shuffle. C is executed, and the previously generated generation matrix is disrupted through the file, namely, the order of the triples is disrupted; and finally, a glove. C training word vector is executed. Before training, you need to make and compile it, and then go through the above steps.

ELMo feature extraction 4.2.4. The process of ELMo feature extraction is as follows: (1) Learn the complex characteristics of word usage through multi-layer LSTM. (2) Through pre-train + Fine tuning, pre-train is performed on a large corpus first, and then fine tuning is performed on the corpus of the downstream task, so that the changes of these complex usages in different contexts can be learned.

BERT Feature Extraction 4.2.5. The initial word vector of each word in the text is used as the input of the BERT model. The vector can be pre-trained by word embedding model as an initial value, and can also be randomly initialized.Each character or word in the text is converted into a one-dimensional initial word vector by querying the word vector through the vocabulary, and the text vector (the value of the vector is automatically learned in the model training process, which is used to characterize the global semantic information of the text and is integrated with the semantic information of a single character/word) and the position vector are extracted. Add the three vectors according to Attention weight as the input of the model, introduce Msaked LM, use the two-way language model for pre-training, and then solve the downstream task through the Fine-tuning mode. The output of BERT model is the vector representation matrix of each character or word, which integrates the semantic information of the full text.

Comparison results of word embedding model 4.3. The scientific literature and the corresponding novelty search formula needed in the process of generating search formula based on antagonistic learning are preprocessed, and the structure, advantages and disadvantages of the word vector embedding model are compared and summarized in Table 3. Through the comparative study, it is found that the automatic construction of novelty search formula through BERT model can solve the problem of polysemy of novelty search formula. The semantic similarity is calculated by cosine similarity, which makes the novelty retrieval more accurate. In the course of the study, it is found that the Bert model also has defects. The vectors generated by the Bert model are dynamic and not interpretable, resulting in the results can not be

reproduced.

Model	Model structure	Information	Advantages	shortcomings
One- hot	N-bit status register	One-bit effective coding	The coding method is intuitive, and each element corresponds to a feature, which is convenient for calculating the loss function and accuracy.	It wastes memory, is not conducive to computation, has sparse features, and does not consider the semantic relationship of context.
Word2 vec	CBOW/ skip-gram	The target word is input and mapped to the hidden layer to predict the probability of the word.	Fixed size left and right window words can be mapped to dense vector	The semantic information of the words outside the window is lost.
GloVe	Optimization of Word2vec	Word2vec model with integrated information	The co-occurrence matrix is used to consider the fusion of local and global information.	Does not address the existing phenomenon of polysemy.
ELMo	Multilayer LSTM	The pre-trained BiLSTM is used to generate word embedding, which is provided to the downstream model	Semantic and syntactic features are extracted through LSTM to solve polysemy and part of speech correspondence.	The extraction ability of LSTM is weaker than that of Transformer; even BiLSTM only considers unilateral data, so its fusion feature ability is weak, and its parallel computing ability is poor.
BERT	Bidirectional Transformer	Masked LM and Next Sentence Prediction are used to capture word-level and sentence-level representations , respectively	11 NLP tasks were pre-trained and fine -tuned to capture longer range dependencies.	Only 15% of tokens in each batch are predicted, and the convergence speed is slow.

TABLE 3 Comparison of Word Vector Model

5. **Conclusions.** Ten thousand abstracts of dissertations, journal papers and conference papers were downloaded from Wan-fang database, and the training set and

test set were made according to the ratio of 9:1. When applying the retrieval formula generation process to train the retrieval expression generation network, the method of antagonistic learning is used to solve the problem of insufficient training samples. Through the training based on BERT and antagonistic learning model, it is found that the same retrieval formula as the real data can be generated. When the generator and the discriminator in the antagonistic learning model reach Nash equilibrium, the effect is optimal, and the same search formula as the standard set can be output, and when the training times are increased, the number of search formulas generated by the model which are the same as those written by experts does not increase. In this paper, the search terms in the generated search formula are matched through the domain word list and the concept synonym list, and the relevant contents in the generated search formula are recommended for scientific researchers to choose. The following is the search formula generated by the input novelty search point through the BERT model and the antagonistic learning model, as shown in Table 4.

TABLE 4. Partial display of generated results								
Novelty check point	Retrieval							
	Expression							
Purpose To investigate the clinical efficacy of hepatectomy in	Hepatectomy and							
different extent for hepatolithiasis.MethodEighty-six patients	hepatolith and (left							
with intrahepatic bile duct stones were treated with surgical	lateral lobectomy							
operations (left lateral lobectomy, left hemihepatectomy, right	or left							
hemihepatectomy and hepatic segmentectomy), and the	hemihepatectomy							
clinical effects of different hepatectomies were	or segmentectomy)							
compared.Results The residual stone rate of left lateral	and residual rate							
lobectomy was significantly higher than that of left	and (extent of							
hemihepatectomy, right hemihepatectomy and segmentectomy	hepatectomy or							
(P < 0.05 or P < 0.01). The complication rate of right	residual stones							
hemihepatectomy and segmentectomy was higher than that of	after operation or							
left lateral lobectomy and left hemihepatectomy ($P < 0.05$).	segmentectomy)							
The excellent and good rate was 81. 25% (65/80). The								
excellent and good rate of left lateral lobectomy was higher								
than that of left hemihepatectomy, and the difference was								
statistically significant ($P < 0.05$).Conclusion Hepatectomy is								
the most effective method for the treatment of intrahepatic								
bile duct stones. For intrahepatic stones not confined to the								
left lateral lobe, left and right hemihepatectomy and								
segmentectomy are superior to left lateral lobectomy. The								
extent of hepatectomy is closely related to the residual stones								
after operation and the effect of surgical treatment.								
Objective To investigate the clinical value of laparoscopic	Hepatectomy and							
hepatectomy combined with mini-incision hepatectomy.	(liver or							
Methods The clinical data of 13 patients who underwent	haloperidol) and							
laparoscopic hepatectomy combined with small incision	primary liver							
hepatectomy in the Department of Hepatobiliary Surgery of	cancer and liver							

Chongqing Emergency Center from June to December 2017 inflammatory

were retrospectively analyzed. There were 10 cases of primary liver cancer, 2 cases of hepatolithiasis and 1 case of hepatic inflammatory pseudotumor. Results The operation was successfully completed in all the 13 patients, and the operation time was 120270 minutes, with an average of $(176, 92 \pm 51)$. 38) min; The mean blood loss was (146.15 \pm 100.96) ml (50400 ml). One day after the operation, they could get out of bed and resume liquid diet. The mean postoperative hospital stay was (6.30 ± 2.14) days (410 days). There was no death and no operative complication. Conclusion Laparoscopic hepatectomy combined with mini-incision hepatectomy has the advantages of minimal invasion, quick recovery and simple operation, which can be used as an option for clinical hepatectomy and conversion to open surgery during laparoscopic or robotic hepatectomy.

pseudotumor and laparoscopic combined with small incision hepatectomy and (clinical hepatectomy or laparoscopic or robotic hepatectomy)

Purpose To observe the effect of rosiglitazone on premature ovarian failure induced by cyclophosphamide (CTX) in patients with systemic lupus erythematosus (SLE). Methods Eight patients with systemic lupus erythematosus were treated with rosiglitazone and treated with cyclophosphamide, and other secondary factors were excluded. Sex hormone levels (FSH, E2) and disease activity index score (SLEDAI) were detected before treatment, during menstruation, one month and three months after menstruation. Results Menstruation returned to normal in 6 patients after 4 to 6 days of rosiglitazone treatment, in 1 patient after 2 months of continuous rosiglitazone treatment, and in 1 patient without recovery of menstruation. No increase in SLEDAI. Conclusion Rosiglitazone is effective in the treatment of SLE with secondary premature ovarian failure caused by cvclophosphamide, without adverse reactions and adverse effects on the disease.

((Systemic lupus erythematosus or SLE) tic and disorder and (cyclophosphamide or secondary factor rosiglitazone) or and (adverse (reaction or SLEDAI) or menstruation)

Novelty search is very important for novelty search workers. In the era of "big data", the content contained in the Internet is blowout growth, the amount of information is miscellaneous, and the massive data sources greatly increase the burden of scientific and technological workers. Therefore, according to the characteristics of sci-tech novelty retrieval texts, this paper optimizes the process of sci-tech novelty retrieval business, and on this basis, uses the combination of word embedding model and antagonistic learning to generate a comprehensive novelty retrieval formula.

It is found that the method based on antagonistic learning can automatically generate search formulas and provide technical support for novelty retrieval system. It realizes the structured analysis technology of novelty search point text and the automatic generation technology of novelty search formula, so as to help scientific researchers quickly understand the content of the novelty search project and further improve work efficiency. Acknowledgment. This work is partially supported by Mr Liu. The authors also gratefully acknowledge the helpful comments and suggestions of the reviewers, which have improved the presentation.

REFERENCES

- Bengio, Y., Ducharme, R., Vincent, P., & Janvin, C. (2003). A Neural Probabilistic Language Model. The Journal of Machine Learning Research, 2003:1137-1155.
- [2] Collobert, R., & Weston, J. A unified architecture for natural language processing. Proceedings of the 25th International Conference on Machine Learning - ICML 2008, 20(1), 160-167.
- [3] Mikolov, T., Corrado, G., Chen, K., & Dean, J. Efficient Estimation of Word Representations in Vector Space. Proceedings of the International Conference on Learning Representations ICLR 2013:1-12.
- [4] Mikolov, T., Chen, K., Corrado, G., & Dean, J. Distributed Representations of Words and Phrases and their Compositionality. NIPS,2013:1-9.
- [5] Pennington, J., Socher, R., & Manning, C. D. Glove: Global Vectors for Word Representation. Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing.2014:1532–1543.
- [6] Peters M,Neumann M,Iyyer M,et al. Deep contextualized word representations[C]. Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics:Human Language Technologies,Volume1(Long Papers). New Orleans,2018: 2227-2237.
- [7] Devlin J,Chang M W,Lee K,et al. BERT: pre-training of deep bidirectional transformers for language understanding[C]. Proceedings of the 2019 Conference of the Association for Computational Linguistics. Stoudsburg PA:Annual Meeting of the Association for Computational Linguistics, 2019: 4171-4186.
- [8] Lecun Y, Bottou L,Bengio Y,et al. Gradient-based learning applied to document recognition [J], Proceedings of the IEEE,1998, 86(11):2278-2324.
- [9] Eiman, J.L. Finding structure in time[J]. Cognitive science, 1990, 14(2), 179-211.
- [10] Salakhutdinov R. Learning deep generative models[J]. Annual Review of Statistics and Its Application,2015(2):361-385.
- [11] Bengio S,Vinyals O,Jaitly N,et al. Scheduled sampling for sequence prediction with recurrent neural networks[J]. Computer Science,2015(2):1171-1179.
- [12] Huszar F. How (not) to train your generative model:scheduled sampling,likelihood adversary?[J]. Computer Science,2015,7(1):11-18.
- [13] Papineni K,Roukos S,Ward T,et al. BLEU:a method for automatic evaluation of machine translation[C]. Proceedings of the 40th Meeting on Association for Computational Linguistics. Stoudsburg PA:Annual Meeting of the Association for Computational Linguistics,2002:311-318.