

Protein-Protein Interaction Extraction Based on Combinational Learning and Active Learning

21209163

20156

Dalian University of Technology

Protein-Protein Interaction, PPI

PPI

Q-statistic

AIMed 71% F 92.9% AUC

PPI

5

LLL

4

AUC

AIMed

BioInfer

Protein-Protein Interaction Extraction Based on Combinational Learning and Active Learning

Abstract

The continuous development of life science and technology has resulted in an explosive increase of the biomedical literature. Therefore, intelligent methods that can automatically extract useful information from massive literatures are urgently needed. As the Internet is becoming increasingly perfecting, information extraction technology has rapidly developed. And, it has a significant impact on biomedical research. Protein-Protein Interaction (PPI) extraction is an important application of information extraction technology in biomedical field. It is aiming at the potential knowledge mining from the molecular level. Centering on the PPI extraction problem, this paper carries out the following researches.

To solve the features deficiency and the problem of limited decision-making ability brought by a single classifier, this paper proposes a combinational learning method. This approach is focusing on the features design and multi-classifier integration. In feature selection, depending on the sentence contexts and the syntactic structures, rich features are extracted to construct feature vectors besides, the information gain method is used to screen the feature ; in classifier integration, three classifiers with different decision-making schemes and higher precision are chosen and used separately, including Support Vector Machine, Maximum Entropy and Naive Bayes. Then, we integrate the classification results by adopting linear weighted method to make sure that the classifier which performed better has a larger weight. The combinational learning approach gained 71% of the F-score and 92.9% of the AUC-score on AIMed corpus.

The combinational learning approach is only fit for the situation with enough labeled corpus. However, there is relatively few labeled corpus in practice. Thus, in order to solve this problem, this paper proposes an active learning on the basis of combinational learning. This method is based on the selection approach of the uncertain samples. It repeatedly selects the samples with most useful information from large amounts of unlabeled corpus and ignores the useless ones. Then, annotate the selected samples and add them to the original training set. The active learning method can not only achieve a better PPI extraction performance, but also reduce the hand-annotated work. The experiments of active learning also achieve higher AUC-scores on the other four corpora except for LLL this approach shows a better generalization performance in large corpora as AIMed and BioInfer.

Key Words Protein-Protein Interaction Extraction; Rich Feature Extraction; Feature Selection; Combinational Learning; Active Learning

	I
Abstract	II
1	1
1.1	1
1.2	2
1.2.1	2
1.2.2	2
1.2.3	3
1.3	4
1.4	5
2	6
2.1	6
2.2	7
2.2.1	8
2.2.2	10
2.3	11
2.4	11
2.5	13
2.6	14
2.6.1	14
2.6.2	15
3	17
3.1	18
3.2	19
3.2.1	20
3.2.2	21
3.3	23
3.4	24
3.5	26
3.5.1	AIMed 26

	3.5.2	AIMed	29
	3.5.3	AIMed	30
	3.6		31
4			32
	4.1		33
	4.1.1		34
	4.1.2		35
	4.1.3		35
	4.2		36
	4.2.1		36
	4.2.2		36
	4.2.3		38
	4.3		38
	4.3.1	AIMed	38
	4.3.2	5	40
	4.4		44
			45
			46
			50
			51
			52

1

1.1

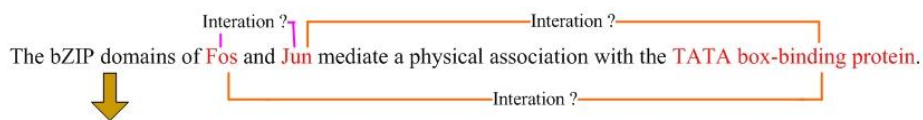
Information extraction, IE

- Protein- Protein Interaction, PPI

PPI

" The bZIP domains of Fos and Jun mediate a physical association with the TATA box-binding protein."

PPI 1.1



- Example1: The bZIP domains of PROT1 and PROT2 mediate a physical association with the TATA box-binding protein. False
- Example2: The bZIP domains of PROT1 and Jun mediate a physical association with the PROT2. True
- Example3: The bZIP domains of Fos and PROT1 mediate a physical association with the PROT2. True

1.1

Fig. 1.1 Example of protein-protein interaction extraction

1.1 3 Fos, Jun, TATA box-binding protein 3

3

PROT1 PROT2 1.1

Example 1

Fos Jun

False Example

2 Example 3

True

PPI

PPI

PPI

AIMed^[1] BioInfer^[2] IEPA^[3] HPRD50^[4] LLL^[5]

PPI

1.2

PPI

PPI

1. 21

[6]

1. 2 2

Part-of-speech Tagging

[7]

1. 2 3

[8]

PPI

PPI

PPI

1

[9]

"

"

Niu

[10]

Liu

[11]

AIMed

53.5%

54.%

F

2

[12]

Walk-weighted Subsequence Kernel
Tree Kernel

Hash Graph Kernel

All-paths Graph Kernels
Airola [13]

PPI

Kim

[14]

PPI

AIMed

F

56.6%

Qian

[15]

PPI

5

Zhang

[16]

AIMed

F

AUC

60.2%

85.3%

Li

[17]

5

PPI

Miwa [18]

PPI

Yang [19]

AIMed F AUC 64.41%

88.46% Li [20]

AIMed F

AUC 69.40% 92.00% Li [21]

AIMed F

4

[22]

F 61.2%

AIMed Qian [23]

AIMed 63.1% F

1.3

PPI

1

PPI

Q-statistic

2

GDep

Stanford Parser Enju

1.4

5

1

PPI

2

3

5

4

PPI

PPI

2

2.1

[24]

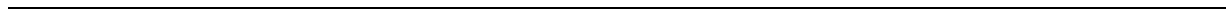
Information Retrieval, IR

[25]

1

2

3



1

2

3

2.2

2.2.1

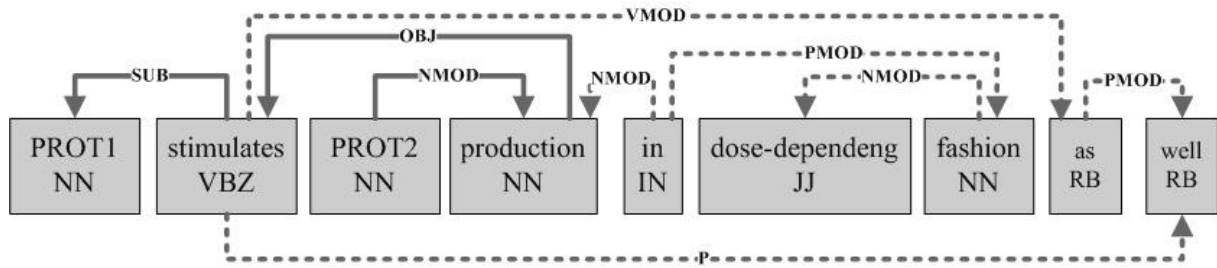
[26]

GDep^[27] GDep Stanford Parser Tsujii

LR

2.1 GDep " PROT1 stimulates PROT2 production in
dose-dependent fashion as well."
" stimulates" 2.1

" stimulates" " VMOD" " stimulates" " as"
" stimulates" " as" " stimulates" " as"



2.1 GDep

Fig. 2.1 Example of the output produced by Gdep

Stanford Parser

[28]

Stanford Parser

Stanford Parser

POS

2.2 2.3 Stanford Parser " PROT1 stimulates PROT2 production
in dose-dependent fashion as well."

2.2

2.3

```
(ROOT
(NP
(NP (CD PROT1) (NNS stimulates))
(NP
(NP (JJ PROT2) (NN production))
(PP (IN in)
(NP
(NP (JJ dose-dependent) (NN fashion))
(CONJP (RB as) (RB well))
(NP (NNP #))))))
```

2.2 Stanford Parser

Fig. 2.2 The phrase structure produced by Stanford Parser


```

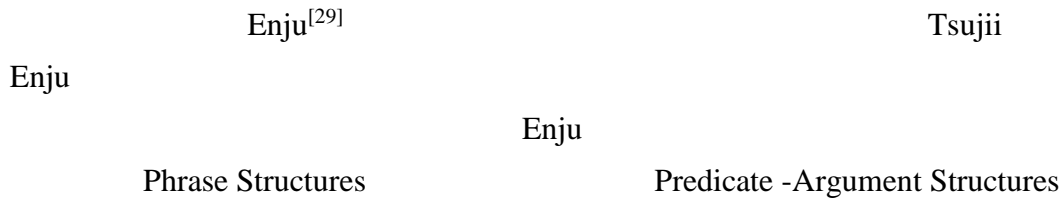
num(stimulates-2, PROT1-1)
root(ROOT-0, stimulates-2)
amod(production-4, PROT2-3)
dep(stimulates-2, production-4)
amod(fashion-7, dose-dependent-6)
prep_in(production-4, fashion-7)
prep_in(production-4, #-10)
conj_and(fashion-7, #-10)

```

2.3 Stanford Parser

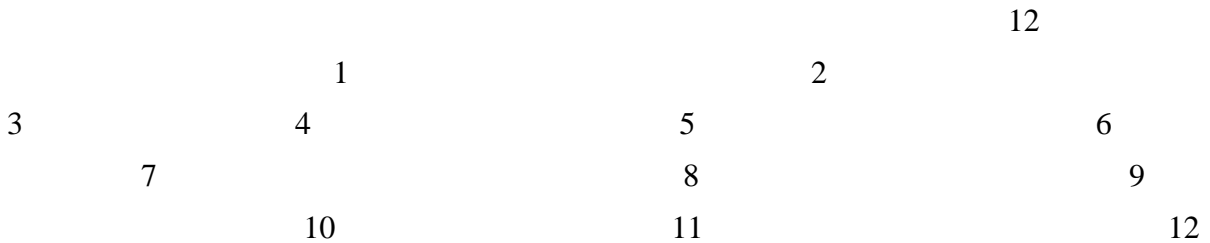
Fig. 2.3 The dependency representation produced by Stanford Parser

2.2.2



2.4

" PROT1 stimulates PROT2 production in dose-dependent fashion as well."



```

ROOT ROOT ROOT ROOT -1 ROOT ROOT stimulates stimulate VBZ VB 1
stimulates stimulate VBZ VB 1 verb_arg12 ARG1 PROT1 prot-NUMBER- CD CD 0
stimulates stimulate VBZ VB 1 verb_arg12 ARG2 production production NN NN 3
PROT2 prot-NUMBER- CD CD 2 adj_arg1 ARG1 production production NN NN 3
dose-dependent dose-dependent JJ JJ 5 adj_arg1 ARG1 fashion fashion NN NN 6
in in IN IN 4 prep_arg12 ARG1 production production NN NN 3
in in IN IN 4 prep_arg12 ARG2 fashion fashion NN NN 6
well well RB RB 8 adj_arg1 ARG1 stimulates stimulate VBZ VB 1
as as RB RB 7 adj_arg1 ARG1 well well RB RB 8
ROOT ROOT ROOT ROOT -1 ROOT ROOT # -SHARP- # -SHARP- 9
# -SHARP- # -SHARP- 9 noun_arg1 ARG1 UNKNOWN UNKNOWN UNKNOWN UNKNOWN -1

```

2.4 Enju

Fig. 2.4 The predicate argument structure produced by Enju

2.3

[30]

1

IG ^[31]

2

(DF)^[32]

t,
DF

t

DF

DF

2.4

[33]

1

SVM 20 90 Vapnik
 SVM VC
 SVM SVM
 SVM SVM
 SVM SVM

[36]

2.1 ~ 2.4

Linear Function

$$K(x, y) = \langle x, y \rangle \quad (2.1)$$

Polynomial Function

$$K(x, y) = \langle x, y \rangle^d \quad (2.2)$$

Radial Basis Function, RBF

$$K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right) \quad (2.3)$$

Sigmoid

Sigmoid Function

$$K(x, y) = \tanh(\beta \langle x, y \rangle) \quad (2.4)$$

2

Maximum Entropy, ME

E.T.Jaynes 1957

[37]

3

[38]

[35]

2.5

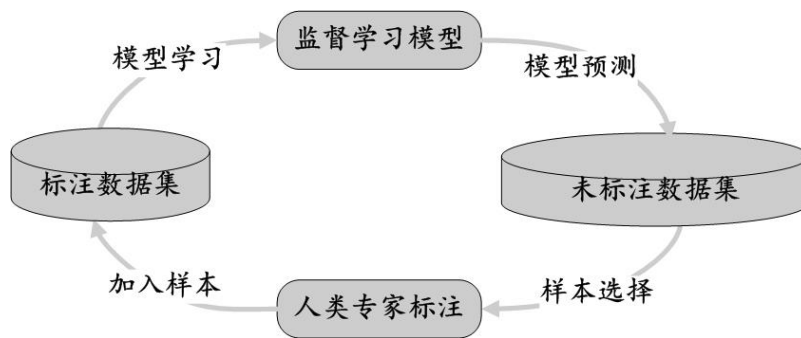
Active Learning

[39]

2.5

[40]

" "



2.5

Fig. 2.5 Framework of active learning

2.6

2.6.1

F-score, F AUC AUC-score, AUC P Precision R Recall F

0 1

1

2.5

2.6

$$P = \frac{TP}{TP + FP} \quad (2.5)$$

$$R = \frac{TP}{TP + FN} \quad (2.6)$$

TP True Positive ; FP False
 Positive ; FN False Negative

$$F = \frac{2 \cdot P \cdot R}{P + R} \quad (2.7)$$

AUC Area Under the Curve
^[41] AUC ROC Receiver Operating Characteristic Curve
 AUC 0.5 1 1
 AUC 2.8

$$AUC = \frac{\sum_{i=1}^n r_i \cdot \frac{n_+ - h_i + 1}{n_+ + n_-}}{n_+ \cdot n_-} \quad (2.8)$$

n_+ n_- r_i i

2.6.2

AIMed BioInfer HPRD50 IEPA LLL

5

XML

True False 5

AIMed 225 PubMed DIP

	200					
	BiInfer					
BiInfer		5				
	HPRD50		HRPD	Human Protein Reference Database		50
				ProMiner		
IEPA	200	PubMed				
LLL	LLL05	Learning Language in Logic 2005				
	2.1			5	PPI	
		5			AIMed	
	5			PPI		
HPRD50			8921	LLL	27	5
	LLL			1	1	AIMed

2.1 5 PPI
 Tab. 2.1 Statistics of the five PPI corpora

AIMed	225	1955	5834	1000	4834
BiInfer	836	1100	8921	2437	6484
HPRD50	50	145	433	163	270
IEPA	200	486	817	335	482
LLL	45	77	330	164	166

3

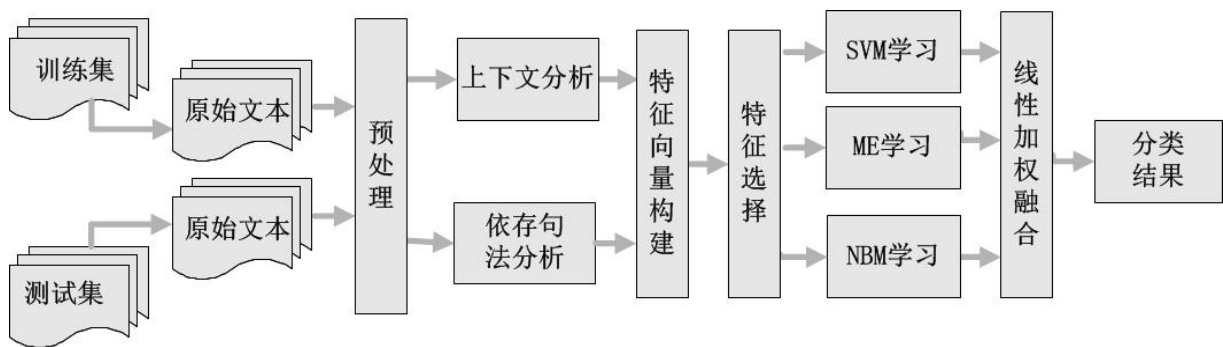
PPI

PPI

PPI

Q-statistic

3.1



3.1

Fig. 3.1 Framework of PPI extraction based on combination of learners

3.1

PPI

5

XML

a the

" . " - " " " < " % "



$n \quad n \dagger 2$

C_n^2

" PROT1" " PROT2"

3.2

[42]

PPI

Tikk

PPI

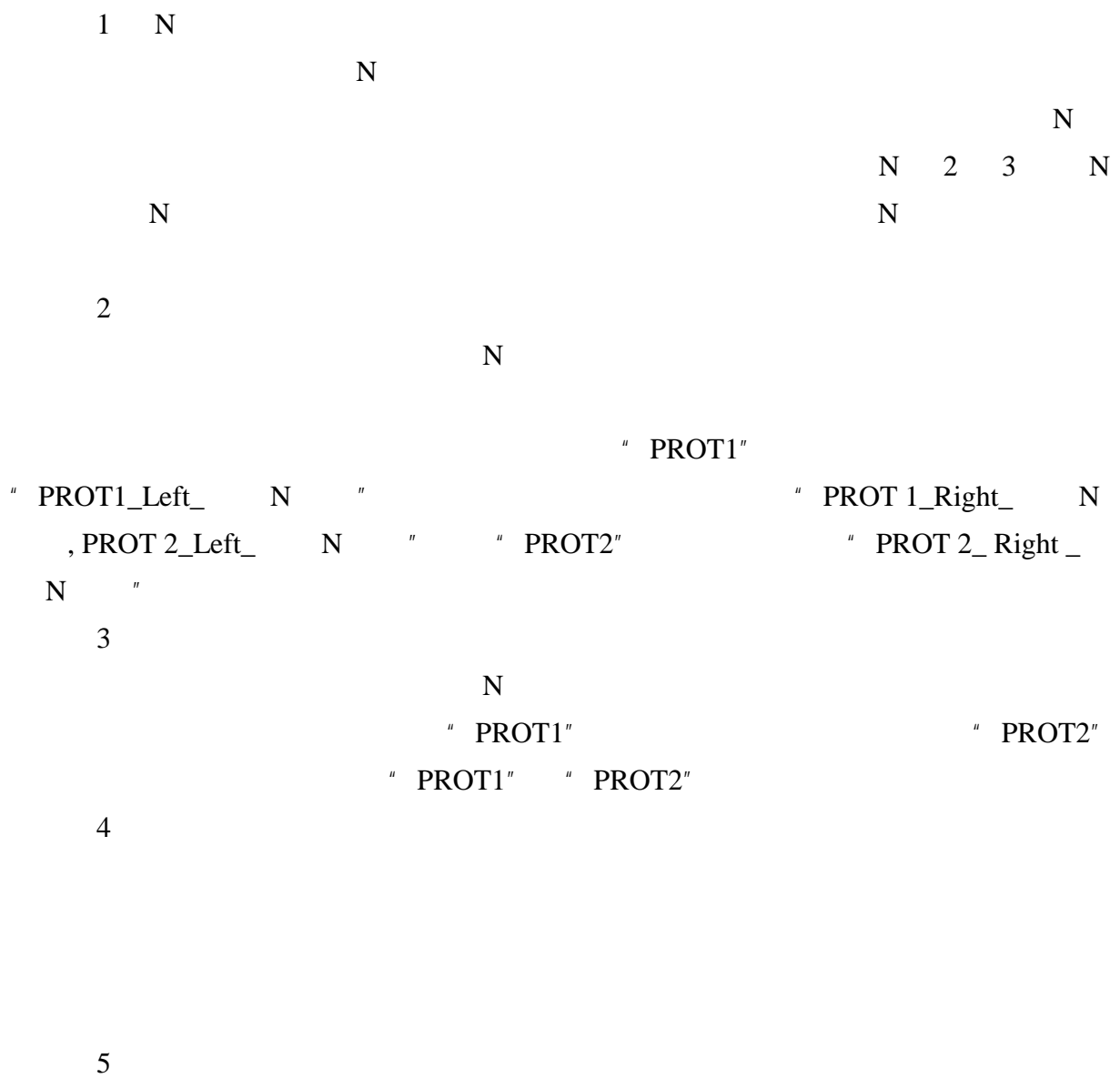
GDep Stanford Parser Enju

GDep

Stanford Parser

Enju

3.2.1



6

5

5

5

7

" neither, hardly, not, nor, never"

3.2.2

1

2

3

4

Stanford Parser

Stanford Parser

5

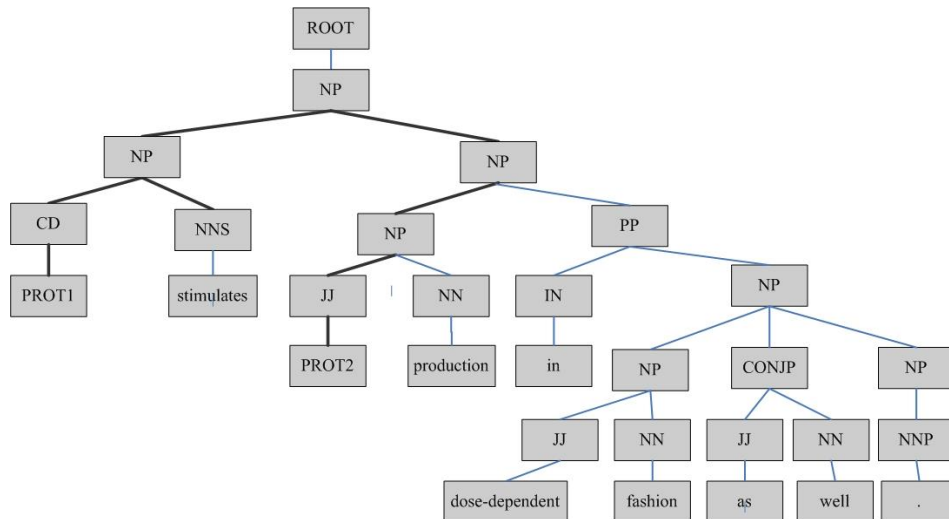
3.2

" PROT1 stimulates PROT2 production in dose-dependent fashion as well."

Stanford

Parser

" PROT1" " PROT2"
 " PROT1" " PROT2"
 (1),(2),(3),(4),(5) Stanford Parser
 Stanford Parser



3.2

Fig. 3.2 Structure path between two proteins

6

" PROT1" " PROT2"

GDep

PROT1 PROT2

Brown

VB VBZ VBD

0

2.1

PROT1 PROT2

1

7 walk

			" PROT1 "	" PROT2 "		
	"	...	"			
				Kim	[14]	e-walk
v-walk				"	" "	"
	(6)	(7)	GDep			GDep
	8					
Zhang	[43]			Shortest Path-enclosed Tree, SPT		5
				Enju		
	"	"	"	"		...
9						
	(8)	(9)	Enju			Enju

3.3

$$\begin{aligned}
 & 3.1 \\
 & IG f_{i|t} = 1 - H f_{i|t} = - \sum_{i=1}^m p f_{i|t} \log p f_{i|t} \quad \checkmark \\
 & p f_{i|t} = \sum_{i=1}^m p f_{i|t} \log p f_{i|t} \checkmark \quad p f_{i|t} = - \sum_{i=1}^m p f_{i|t} \log p f_{i|t} \\
 & p f_{i|t} \quad c_i \quad p f_{i|t} \quad t \\
 & p f_{i|t} \quad t \quad c_i \quad p f_{i|t}
 \end{aligned} \tag{3.1}$$

$$t \quad p_{c_i|t} \quad t \quad c_i \quad m$$

3.4

PPI

0.5

Q-statistic^[44]

$Q_i \quad Q_j$

$$Q_{ij} = \frac{p^{11}(t, j) p^{00}(t, j) - p^{01}(t, j) p^{10}(t, j)}{p^{11}(t, j) p^{00}(t, j) + p^{01}(t, j) p^{10}(t, j)} \quad (3.2)$$

$$\begin{array}{rcc}
 p^{11}(t, j) & & Q_i \quad Q_j \\
 p^{10}(t, j) & & Q_i \quad Q_j \\
 p^{01}(t, j) & & Q_i \quad Q_j \\
 p^{00}(t, j) & & Q_i \quad Q_j \\
 Q_{ij} & -1 \quad 1 & Q_{ij} = 0
 \end{array}$$

L

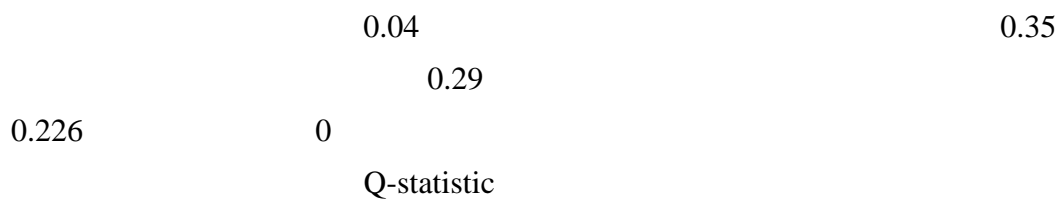
Q_{ij}

3.3

$$\bar{Q} = \frac{2}{L(L-1)} \sum_{i=1}^{L-1} \sum_{j=1}^L Q_{ij} \quad (3.3)$$

SVM

3.2 3.3



$$\sum_{n=1}^N K(x, x') = \frac{1}{n} \sum_{n=1}^N g_n K_n(x, x') \quad (3.4)$$

$\sum_{n=1}^N g_n = 1, \quad g_n \geq 0, \quad \sum_{n=1}^N g_n = 1$

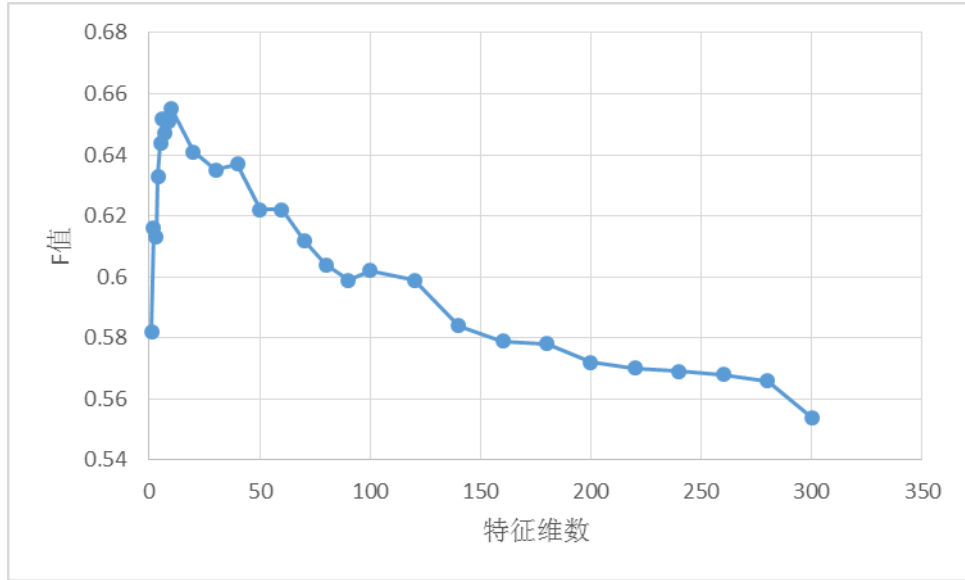
3.1 AI Med
 Tab. 3.1 Experiment results of different dimension features on AIMed

	($\times 1000$)	P	R	F	AUC
1	1	0.587	0.581	0.582	0.866
2	2	0.599	0.642	0.616	0.887
3	3	0.587	0.647	0.613	0.891
4	4	0.592	0.689	0.633	0.901
5	5	0.596	0.708	0.644	0.904
6	6	0.600	0.723	0.652	0.905
7	7	0.591	0.721	0.647	0.906
8	8	0.588	0.741	0.652	0.906
9	9	0.586	0.741	0.651	0.907
10	10	0.586	0.749	0.655	0.908
11	20	0.555	0.766	0.641	0.908
12	30	0.542	0.773	0.635	0.906
13	40	0.547	0.771	0.637	0.906
14	50	0.525	0.771	0.622	0.902
15	60	0.519	0.785	0.622	0.900
16	70	0.509	0.780	0.612	0.897
17	80	0.498	0.778	0.604	0.893
18	90	0.491	0.781	0.599	0.892
19	100	0.493	0.789	0.602	0.887
20	120	0.485	0.790	0.599	0.886
21	140	0.473	0.780	0.584	0.880
22	160	0.465	0.783	0.579	0.878
23	180	0.462	0.788	0.578	0.876
24	200	0.458	0.781	0.572	0.870
25	220	0.454	0.780	0.570	0.870
26	240	0.453	0.781	0.569	0.869
27	260	0.451	0.781	0.568	0.868
28	280	0.449	0.784	0.566	0.868
29	300	0.436	0.775	0.554	0.857

10000

300000

10000

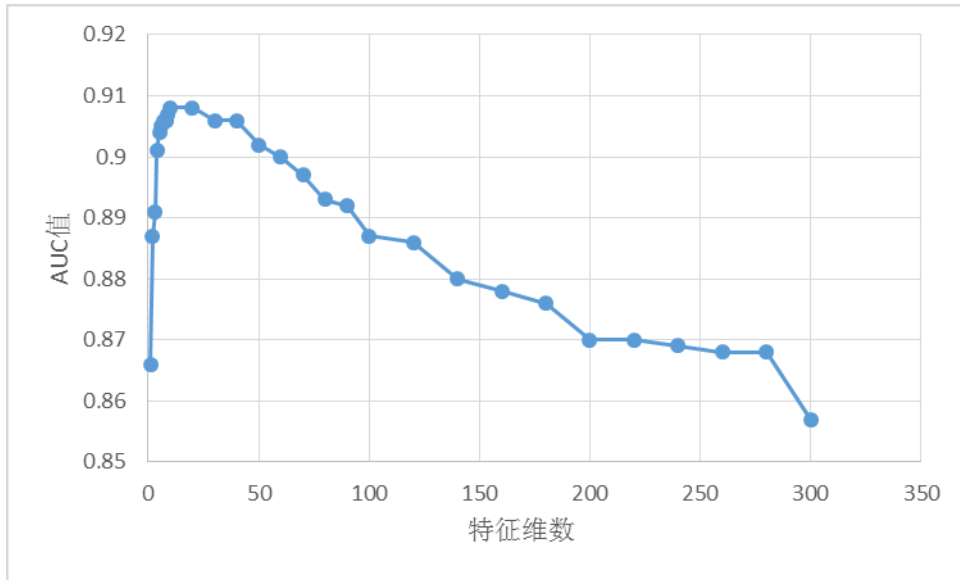


3.3

AIMed

F

Fig. 3.3 Results of F-score between different dimension features on AIMed



3.4

AIMed

AUC

Fig. 3.4 Results of AUC-score between different dimension features on AIMed

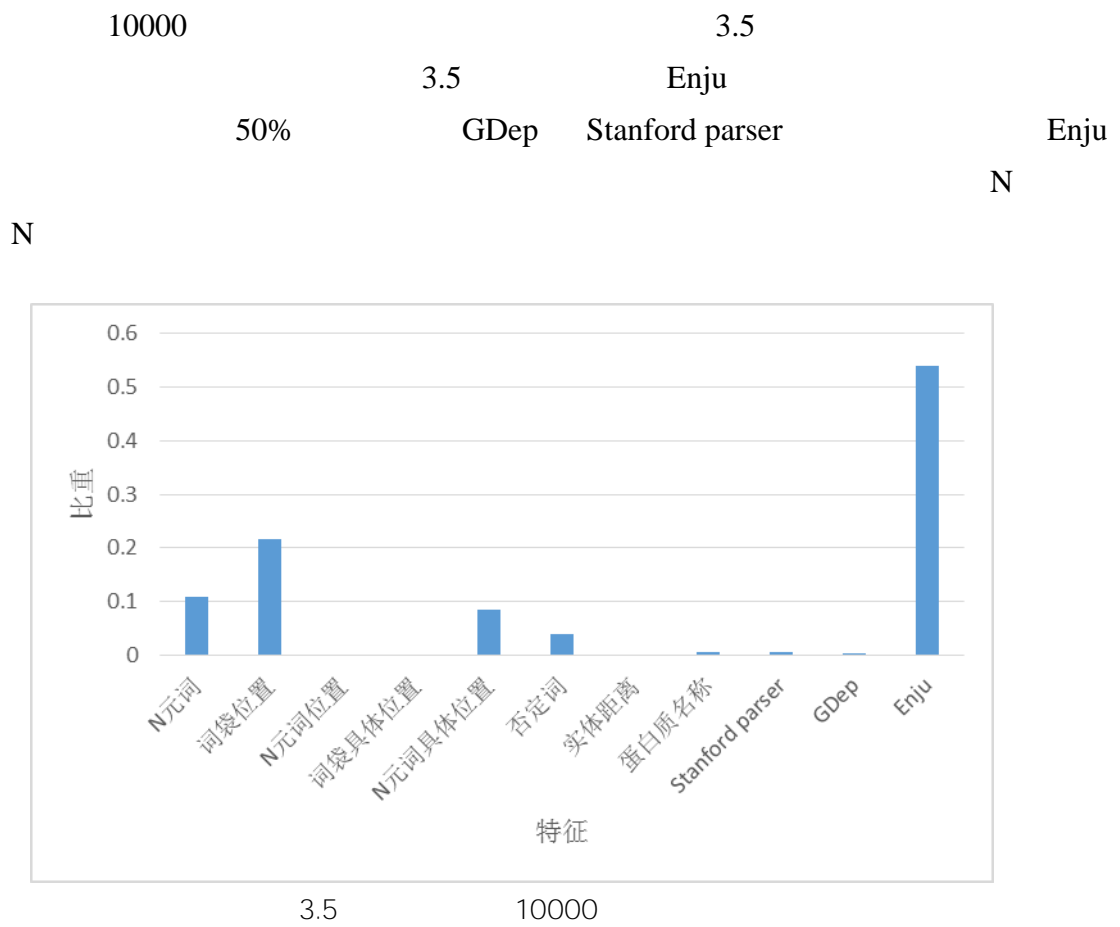
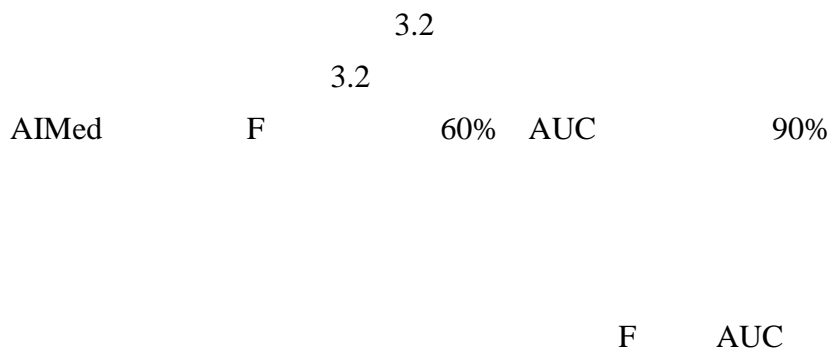


Fig. 3.5 The proportion of different features 10,000-dimensional feature

3.5.2

AIMed



3.6

PPI

Q-statistic

AIMed	F	AUC	71%	92.9%		
PPI	5				LLL	4
	AUC				PPI	

4

3

PPI

PPI

PPI

PPI

SVM

4.1

PPI

PPI

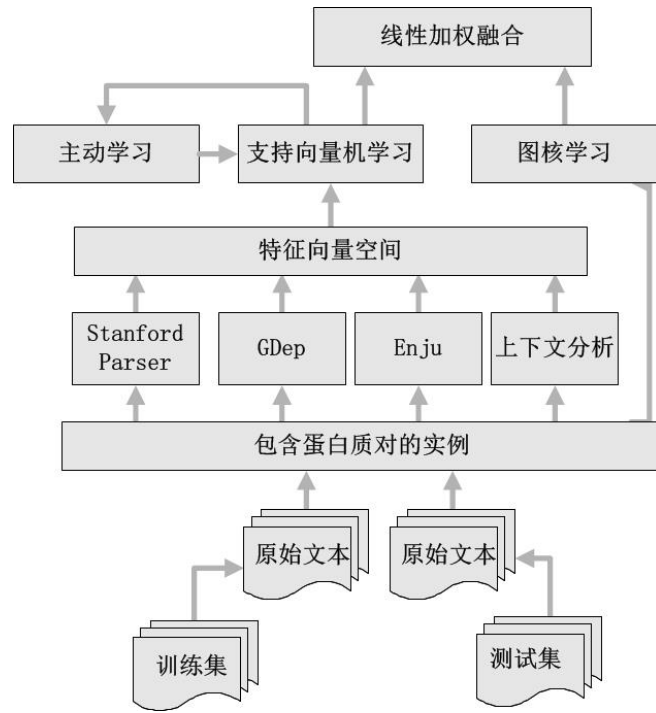
PPI

1

2

3

3.3



4.1

Fig. 4.1 Framework of PPI extraction based on active learning

4.1

PPI

PPI

4.1.1

[46]

Logistic

5

4

SVM

0.1

4.2

基于不确定样本选择的主动学习算法

Input:

从已标注语料中划分出来的训练集: O_1 ;
剩余的四个语料作为未标注语料集: O_2 ;
从已标注语料中划分出来的测试集: T ;
使用特征向量的SVM分类器模型: S ;
使用特征向量的ME分类器模型: M ;
每次选择的新样本的数量: N 。

Output:

用模型 S 训练 T 得到的准确率 P , 召回率 R , F值以及AUC值。

```
1 训练 $O_1$ , 得到模型 $S$ 和 $M$ 。  
2  Repeat  
3     使用最大熵模型 $M$ 对测试集 $O_2$ 进行分类, 并从测试集 $O_2$ 中选择出  
   模型 $M$ 最不确定的 $N$ 个新样本。  
4     对选出来的 $N$ 个新样本进行标注, 并添加到 $O_1$ 中。  
5     将此 $N$ 个样本从 $O_2$ 中去除掉。  
6     重新训练 $O_1$ , 得到新的分类器模型 $S$ 和 $M$ 。  
7     使用当前的模型 $S$ 预测测试集 $T$ 。  
8     计算 $P$ ,  $R$ , F值以及AUC值。  
9  Until 没有模型 $M$ 最不确定的样本。  
10 Print 最终的准确率 $P$ , 召回率 $R$ , F值以及AUC值。
```

4.2

Fig. 4.2 Active learning algorithm based on sample selection uncertainty

4.1.2

[47]

Committee, QBC Seung^[48] Freund^[49] Query By QBC

4.1.3

[50]

4.2

PPI

4.2.1

PPI

PPI

PPI

1

2

PPI

svm^{light}[51]

4.2.2

Airor

[13]

PSS

LOS

PSS

IP

IP

PSS

LOS

4.2.3

PPI

3

3.1

0.7 0.3

0 1

0.1

4.3

4.3.1 AIMed

AIMed

PPI

PPI

AIMed

Stanford Parser GDep Enju

PPI

Stanford

GDep

Enju

4.2

F AUC

4.2

GDep F AUC 8.5 5.5

PPI Stanford F

Enju F 65.3% AUC 90.3%

4.2 AIMed

Tab. 4.2 The experiment results after different procedures on AIMed corpora

	P	R	F	AUC
	0.559	0.535	0.541	0.829
+GDep	0.622	0.636	0.626	0.884
+Stanford	0.618	0.652	0.632	0.885
+Enju	0.649	0.665	0.653	0.903
+	0.647	0.688	0.662	0.905
+	0.673	0.663	0.664	0.909

GDep Enju

F 10.6 AUC 7.3

PPI

AIMed

F /66.4% AUC /90.9%

4.3 PPI AIMed

F AUC

4.3 AIMed F

66.4% AUC 90.9% Miwa [18]

60.8% F 86.8% AUC Kim [14]

F 56.6%

F	AUC	Yang	[19]
4	71.1%	Qian	[15]
F	63.1%	Qian	[23]

4.3 AIMed

Tab. 4.3 Performance comparison between different methods on AIMed corpora

	P	R	F	AUC
Miwa[18]	0.550	0.688	0.608	0.868
Kim [14]	0.614	0.533	0.566	-
Yang [19]	0.577	0.711	0.644	0.885
Qian [15]	0.591	0.576	0.581	0.833
Qian [23]	0.651	0.613	0.631	-
	0.673	0.663	0.664	0.909

4.3.2

5

5

4.4

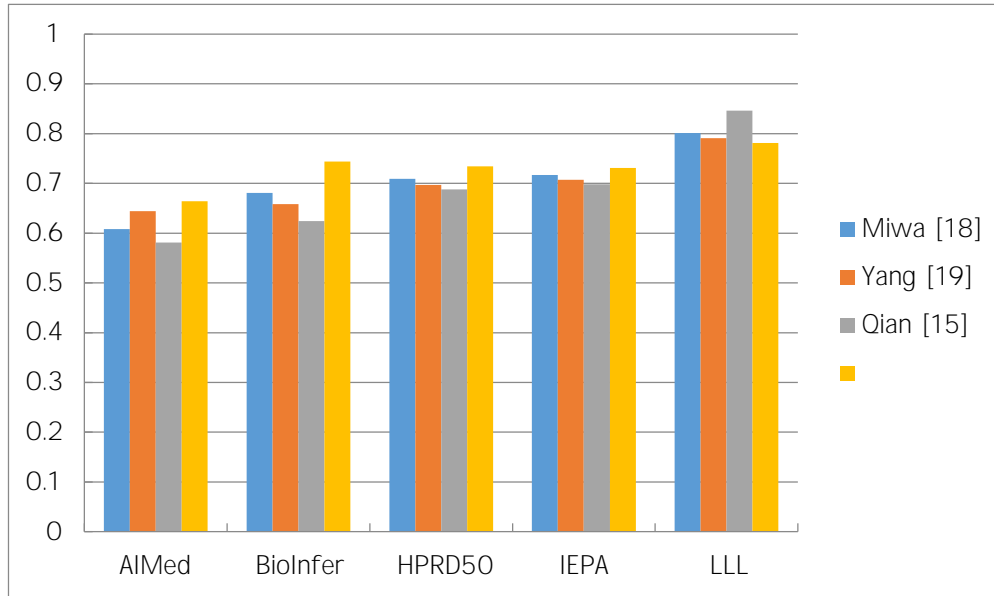
4.4 5

Tab. 4.4 Performance comparison between different methods on five corpora

		AIMed	BiInfer	HPRD50	IEPA	LLL
Miwa[18]	F	0.608	0.681	0.709	0.717	0.801
	AUC	0.868	0.859	0.822	0.844	0.863
Yang [19]	F	0.644	0.658	0.697	0.707	0.791
	AUC	0.885	0.850	0.744	0.757	0.830
Qian [15]	F	0.581	0.624	0.688	0.698	0.846
	AUC	0.833	0.836	0.837	0.828	0.899
	F	0.664	0.744	0.734	0.731	0.781
	AUC	0.909	0.921	0.845	0.855	0.866

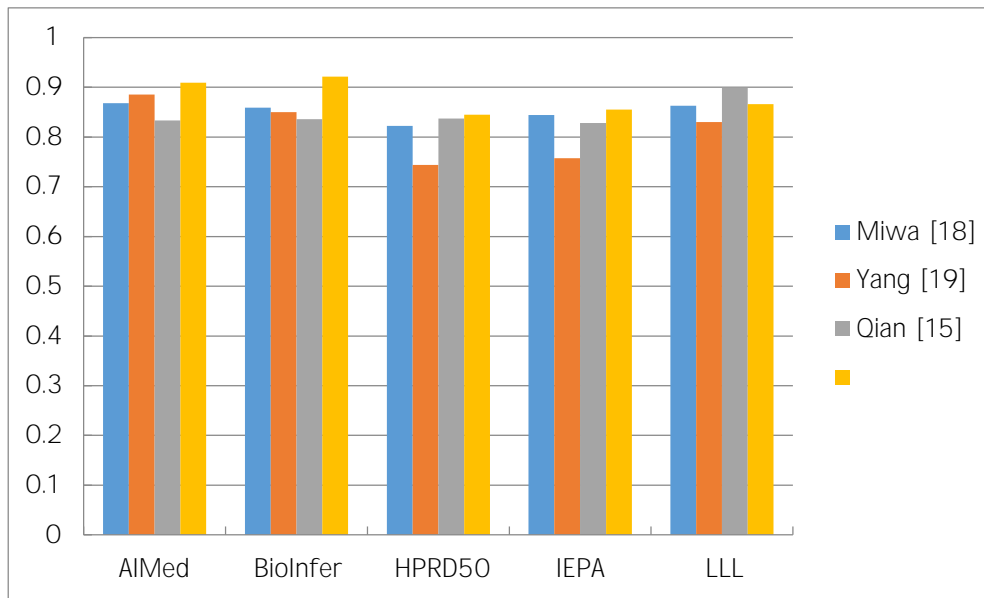
4.4 4.5

F AUC



4.4 F 5

Fig. 4.4 Comparison results of F-score between different methods on five corpora



4.5 AUC 5

Fig. 4.5 Comparison results of AUC-score between different methods on five corpora

4.4 4.5 4 AIMed BioInfer HPRD50
 IEPA F AUC BioInfer
 F AUC
 LLL

5 PPI

[52,53]

5 4
 4.5 F 4.6 AUC
 4.5 F

Tab. 4.5 Cross-corpus results comparison measured with F-score

	AIMed		BioInfer		HPRD50		IEPA		LLL		Avg.rank
	F	Rank	F	Rank	F	Rank	F	Rank	F	Rank	
AIMed	-	-	0.530	2	0.674	2	0.672	2	0.734	2	2
BioInfer	0.485	1	-	-	0.685	1	0.740	1	0.753	1	1
HPRD50	0.425	2	0.471	4	-	-	0.640	4	0.706	4	3.5
IEPA	0.385	3	0.548	1	0.643	4	-	-	0.708	3	2.75
LLL	0.345	4	0.491	3	0.650	3	0.646	3	-	-	3.25

4.6 AUC

Tab. 4.6 Cross-corpus results comparison measured with AUC score

	AIMed		BioInfer		HPRD50		IEPA		LLL		Avg.rank
	AUC	Rank	AUC	Rank	AUC	Rank	AUC	Rank	AUC	Rank	
AIMed	-	-	0.702	2	0.804	1	0.758	2	0.782	4	2.25
BioInfer	0.792	1	-	-	0.801	2	0.841	1	0.841	1	1.25
HPRD50	0.714	2	0.652	4	-	-	0.725	3	0.810	3	3
IEPA	0.668	3	0.721	1	0.764	3	-	-	0.827	2	2.25
LLL	0.610	4	0.658	3	0.752	4	0.723	4	-	-	3.75

Rank 4.5

AIMed Rank 1 F LLL Rank

F Rank 4 Avg.rank Rank

Avg.rank PPI

4.5 AIMed F

48.5% BioInfer AIMed F 66.4%

F AUC

BioInfer IEPA

BioInfer IEPA

4.6 AUC 61.0% 84.1%

AUC

AUC 4.4 AUC 84.5% 92.1%

4.5 4.6 F AUC Avg.rank

4.6 4.7 Miwa [13]

Avg.rank 4.6

4.7

Avg.rank

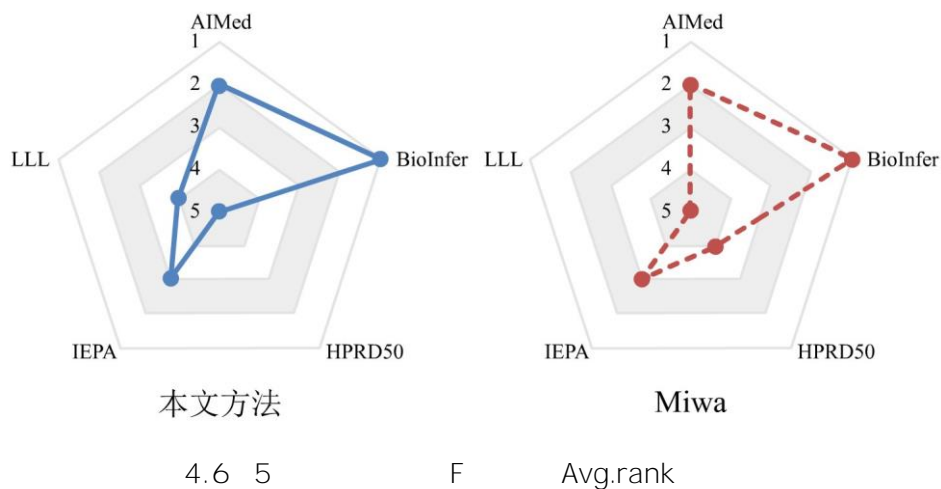
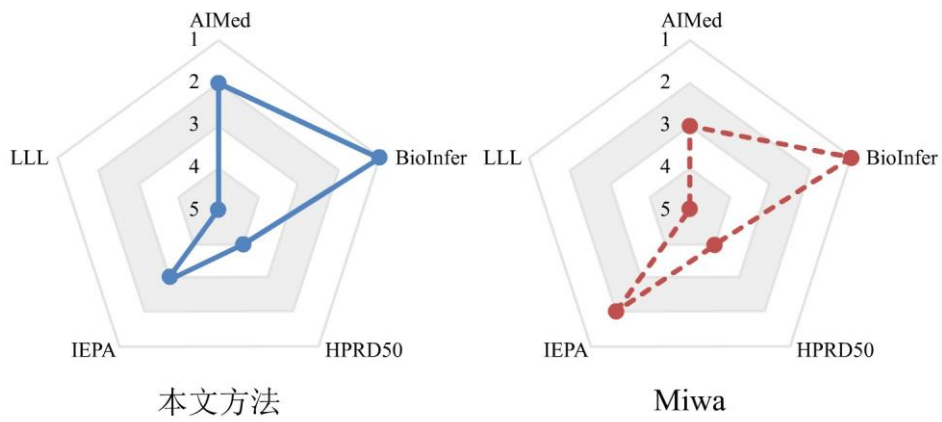


Fig.4.6 Comparison results between different Avg.rank values of F-score on five corpus



4.7 5 AUC Avg.rank
 Fig.4.7 Comparison results between different Avg.rank values of AUC-score on five corpus

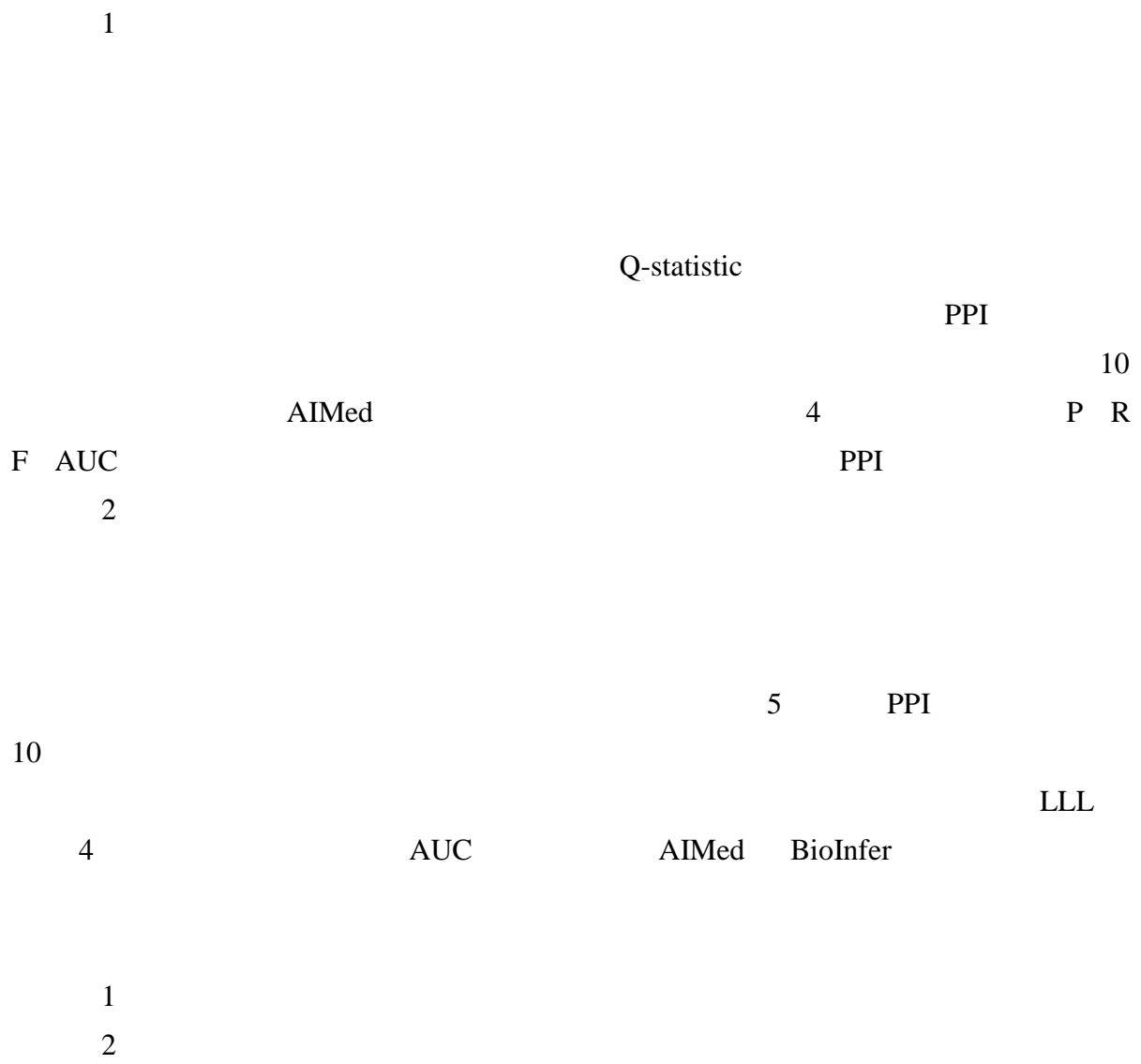
4.6 4.7 Miwa [13]
 BioInfer
 HPRD50 LLL Avg.rank
 PPI

4.4

Gdep

Enju

PPI



-
- [1] Bunescu R, Ge R, Kate R J, et al. Comparative experiments on learning information extractors for proteins and their interactions [J]. *Artificial intelligence in medicine*, 2005, 33(2): 139-155.
- [2] Pyysalo S, Ginter F, Heimonen J, et al. BioInfer: a corpus for information extraction in the biomedical domain [J]. *BMC bioinformatics*, 2007, 8(1): 50.
- [3] Ding J, Berleant D, Nettleton D, et al. Mining MEDLINE: abstracts, sentences, or phrases? [C]//Pacific Symposium on Biocomputing. 2002, 7: 326-337.
- [4] Fundel K, Küfner R. Relation extraction using dependency parse trees [J]. *Bioinformatics*, 2007, 23(3): 365 -371.
- [5] Nédellec C. *ing language in logic* -genic interaction extraction challenge [C]//Proceedings of the 4th Learning Language in Logic Workshop (LLLO5), 2005, 7.
- [6] Bunescu R, Mooney R, Ramani A, et al. Integrating co -occurrence statistics with information extraction for robust retrieval of protein interactions from Medline [C]//Proceedings of the Workshop on Linking Natural Language Processing and Biology: Towards Deeper Biological Literature Analysis. Association for Computational Linguistics, 2006: 49 -56.
- [7] Huang M, Zhu X, Hao Y, et al. Discovering patterns to extract protein -protein interactions from full biomedical texts [J]. *Bioinformatics*, 2004, 20(18): 36043612.
- [8] Shiguihara -J u á r e z P N , d e A n d r a d e L o p e s A . L e a r n i n g trees for extraction of protein -protein interaction [C]//Proceedings of the 14th international conference on Computational Linguistics and Intelligent Text Processing-Volume 2. Springer -Verlag, 2013: 347-358.
- [9] Mitsumori T, Murata M, Fukuda Y, et al. Extracting protein -protein interaction information from biomedical text with SVM [J]. *IEICE Transactions on Information and Systems*, 2006, 89(8): 24642466.
- [10] Niu Y, Otasek D, Jurisica I. Evaluation of linguistic features useful in extraction of interactions from PubMed; application to annotating known, high -throughput and predicted interactions in I2D [J]. *Bioinformatics*, 2010, 26(1): 111 -119.
- [11] Liu B, Qian L, Wang H, et al. Dependency -driven feature -based learning for extracting protein -protein interactions from biomedical text [C]//Proceedings of the 23rd International Conference on Computational Linguistics: Posters. Association for Computational Linguistics, 2010: 757 -765.

-
- [12] Cristianini N, Shawe - Taylor J. An introduction to support vector machines and other kernel -based learning methods [M]. Cambridge university press, 2000.
- [13] Airola A, Pysalo S. Graph kernel for protein -protein interaction extraction with evaluation of cross -corpus learning [J]. BMC bioinformatics, 2008, 9(Suppl 11): S 2.
- [14] Kim S, Yoon J, Yang J, et al. Walk -weighted subsequence kernels for protein -protein interaction extraction [J]. BMC bioinformatics, 2010, 11(1): 107.
- [15] Qian L, Zhou G. Tree kernel -based protein-protein interaction extraction from biomedical literature [J]. Journal of biomedical informatics, 2012, 45(3):535 - 543.
- [16] Zhang Y, Lin H, Yang Z, et al. Neighborhood hash graph kernel for protein -protein interaction extraction [J]. Journal of biomedical informatics, 2011, 44(6): 10861092.
- [17] Li L, Guo R, Jiang Z, et al. An Approach to Improve Kernel -Based Protein-Protein Interaction Extraction by Learning from Large -Scale Network Data [J]. Methods, 2015.
- [18] Miwa M, Sæt re R. Protein interaction extraction by leveraging multiple kernels and parsers [J]. International journal of medical informatics, 2009, 78(12): e39 -e46.
- [19] Yang Z, Tang N, Zhang X, et al. Multiple kernel learning in protein -protein interaction extraction from biomedical literature [J]. Artificial intelligence in medicine, 2011, 51(3): 163 -173.
- [20] Li L, Zhang P, Zheng T, et al. Integrating Semantic Information into Multiple Kernels for Protein -Protein Interaction Extraction from Biomedical Literatures [J]. PloS one, 2014, 9(3): e91898.
- [21] Li L, Jin L, Zheng J, et al. The Protein-Protein Interaction extraction based on full texts [C]//Bioinformatics and Biomedicine (BIBM), 2014 IEEE International Conference on. IEEE, 2014: 493 496.
- [22] Bui Q C, Katrenko S, Sloot P M A. A hybrid approach to extract protein - protein interactions [J]. Bioinformatics, 2011, 27(2): 259 -265.
- [23] Qian W, Fu C, Cheng H. Semi -supervised method for Extraction of Protein -Protein Interactions using hybrid model [C]//Proceedings of the 2013 Third International Conference on Intelligent System Design and Engineering Applications. IEEE Computer Society, 2013: 12681271.
- [24] , , . [J]. , 2009, 44(21): 166170.
- [25] . [D]. : ,2008.
- [26] , . [J]. ,2009,11(2):100-112.

-
- [27] Miva M, Pyysalo S, Hara T, et al. A comparative study of syntactic parsers for event extraction [C]//Proceedings of the 2010 Workshop on Biomedical Natural Language Processing. Association for Computational Linguistics, 2010: 37 -45.
- [28] Stanford Berkeley [J]. ,2013,9(8):1985-1986.
- [29] Miya o Y , Sa ga e E. Evaluating contributions of natural language parsers to protein -protein interaction extraction [J]. Bioinformatics, 2009, 25(3): 394-400.
- [30] [J]. , 2012, 27(2): 161-166.
- [31] [J]. , 2012, 48(27): 119-122.
- [32] [D]. , 2010.
- [33] [D]. : , 2009.
- [34] Bartlett P L, Traskin M. Adaboost is consistent [J]. Journal of Machine Learning Research, 2007, 8: 2347-2368.
- [35] Sun Y, Todorovic S, Li J, et al. Unifying the error -correcting and output -code adaboost within the margin framework [C]//Proceedings of the 22nd international conference on Machine learning. ACM, 2005: 872-879.
- [36] [M]. : 2012.
- [37] Berger AL, Pietra VJD, Pietra SAD. A maximum entropy approach to natural language processing [J]. Computational linguistics, 1996, 22(1): 39 -71.
- [38] Domingos P, Pazzani M. On the optimality of the simple Bayesian classifier under zero-one loss [J]. Machine learning, 1997, 29(2 -3): 103-130.
- [39] [J]. , 2011, 37(8): 954-962.
- [40] [D]. : , 2013.
- [41] [D]. : , 2012.
- [42] Tikk D , Solt I, Thomas P, et al. A detailed error analysis of 13 kernel methods for protein -protein interaction extraction [J]. BMC bioinformatics, 2013, 14(1):12.
- [43] Zhang M, Zhang J, Su J, et al. A composite kernel to extract relations between entities with both flat and structured features [C]//Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2006: 825 -832.

-
- [44] Yule G U. On the association of attributes in statistics: with illustrations from the material of the childhood society, &c[J]. Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character, 1900: 257319.
- [45] Giuliano C, Lavelli A, Romano L. Exploiting shallow linguistic information for relation extraction from biomedical literature [C]//EACL, 2006, 18: 401 - 408.
- [46] , , . [J]. ,200845(): 300304
- [47] , , . [J]. ,2010,46(13): 11-14.
- [48] Seung H S, Opper M, Sompolinsky H. Query by committee [C]//Proceedings of the fifth annual workshop on Computational learning theory. ACM, 1992: 28294.
- [49] Freund Y, Seung H S, Shamir E, et al. Selective sampling using the query by committee algorithm [J]. Machine learning, 1997, 28(2 -3): 133-168.
- [50] , , . [J]. ,2013, 26(12):1121-1129.
- [51] Joachims T. SVM light support vector machine [EB/OL]. 2002. <http://svmlight.joachims.org> .
- [52] Van Landeghem S, Saeys Y, De Baets B, et al. Extracting protein -protein interactions from text using rich feature vectors and feature selection [C]//3rd International symposium on Semantic Mining in Biomedicine (SMBM 2008)Turku Centre for Computer Sciences (TUCS), 2008: 77-84.
- [53] Zhang Y, Lin H, Yang Z, et al. Hash subgraph pairwise kernel for protein -protein interaction extraction [J]. IEEE/ACM Transactions on Computational Biology and Bioinformatics (TC BB), 2012, 9(4): 1190-1202.

1

2015 11

2 WANG J, LIU M, LIN H, et al. Combining Active Learning and Composite Kernel for Protein-protein Interaction Extraction[J]. Journal of Computational Information Systems, 2015, 11(8): 2823-2832.

3

2014SR156358