

## 15. フレーズテーブルの作り方

内山将夫@NICT  
mutiyama@nict.go.jp

## フレーズテーブルの作成手順

1. 初期単語アラインメント
2. 単語アラインメントの修正
3. 単語単位の翻訳確率の推定
4. フレーズ対応の抽出
5. フレーズへのスコアの付与

### 結果

英語句 ||| 日本語句 ||| スコア 1 スコア 2 ...

これらの句のスコアを上手く組み合わせて文全体のスコアを得る .

## 初期単語アラインメント

GIZA++ というフリーソフトウェアを利用することが多い．これを，たとえば，日英の方向に，

- IBM Model-1 (紹介済)
- HMM Model
- IBM Model-3
- IBM Model-4

を順番に適用していくことにより，日英の各文対応について，1対 $n$ の単語対応を得ることができる．  
英日にもすることにより，両方向に，1対 $n$ の単語対応を得る．

=====

単語アラインメントは未解決の問題だが，ある程度まではできるので，とりあえず，そこまでを前提にして，次に進んでいる．

## 日英単語アラインメントの例

	when	the	fluid	pressure	cylinder	31	is	used	,	fluid	is	gradually	applied	.
流体			*											
圧				*										
シリンダ					*									
31						*								
の														
場合	*													
は							*							
流体										*				
が											*			
徐々に												*		
排出												*		
さ													*	
れる													*	
こと													*	
と									*					
なる														
。														*

## 日英単語アラインメントの例

	this	relation	must	be	maintained	even	after	passing	at	least	100,000	sheets	.
そして							*						
、				*									
上記	*												
関係		*											
を								*					
少なくとも										*			
10万											*		
枚											*		
通											*		
紙												*	
し								*					
て						*							
も						*							
維持					*								
し								*					
なけれ			*										
ば			*										
なら			*										
ない			*										
。													*

## 日英単語アラインメントの例

	first	,	a	description	will	be	given	of	the	structure	of	the	wheels	2	.
まず	*														
、		*													
車輪													*		
2														*	
の											*				
構造										*					
について					*										
説明					*										
する					*										
。															*

## 英日単語アラインメントの例

	流体	圧	シリンダ	3 1	の	場合	は	流体	が	徐々に	排出	さ	れる	こと	と	なる
when						*										
the							*									
fluid	*															
pressure		*														
cylinder			*													
31				*												
is													*			
used															*	
,																
fluid								*								
is									*							
gradually										*						
applied											*					
.																

## 英日単語アラインメントの例

	まず	、	車輪	2	の	構造	について	説明	する	。
first	*									
,		*								
a										
description								*		
will								*		
be								*		
given								*		
of					*					
the					*					
structure						*				
of							*			
the									*	
wheels			*							
2				*						
.										*



## 日英英日アラインメントの合成

- アラインメントの積集合 =  $I$   
→ 精度の高い単語対応
- アラインメントの積集合 =  $U$   
→ カバー率の高い単語対応

$I$  を出発点として  $U$  中のアラインメントを加えていく .

## 日英と英日の合成の例

	when	the	fluid	pressure	cylinder	3l	is	used	,	fluid	is	gradually	applied	.
流体			I											
圧				I										
シリンダ					I									
3 l						I								
の														
場合	I													
は		U					U							
流体										I				
が											I			
徐々に												I		
排出												U	U	
さ													U	
れる							U						U	
こと													U	
と								U	U					
なる														
。														I

## 日英と英日の合成の例

	this	relation	must	be	maintained	even	after	passing	at	least	100,000	sheets	.
そして							U						
、				U									
上記	I												
関係		I											
を								U					
少なくとも									U	I			
10万											I		
枚											U	U	
通								U			U		
紙												U	
し							U	U					
て						U							
も						I							
維持					I								
し								U					
なけれ			U										
ば			U										
なら			I	U									
ない			U										
。													I

## 日英と英日の合成の例

	first	,	a	description	will	be	given	of	the	structure	of	the	wheels	2	.
まず	I														
、		I													
車輪													I		
2														I	
の								U	U		U				
構造										I					
について					U						U				
説明				U	I	U	U								
する					U							U			
。															I

# 合成手順

```
#
# アラインメント候補 = {(e_i, f_j)} は e_i + f_j が小さいものから並
# んでいる . そうすることにより , 文頭のものから優先して対応付けてい
# くことができる .
#

日英英日アラインメントの合成 ()
  隣接点集合 = ((-1,0), (0,-1), (1,0), (0,1), (-1,-1), (-1,1), (1,-1), (1,1))
  アラインメント候補 = 積集合 I(日英アラインメント, 英日アラインメント)
  隣接する 1 対 1 の点を追加する ()
  その他の点を加える ()
  アラインメント候補を返す
end

隣接する 1 対 1 の点 N を追加する ()
  while true
    for (j, e)   アラインメント候補
      for (j0, e0)   隣接点 (j,e)   和集合 (日英, 英日)
        if (j0 も e0 もアラインメント候補で使われていない)
          (j0, e0) をアラインメント候補に追加する (そうしても 1 対 1 が壊れない)
        end
      end
    end
  end
  アラインメント候補に追加がなければ終了する
end

#
# なるべく多くの点を加えたいが , 日英共に使われているのは加えない
#
その他の点 0 を加える ()
  for (j,e)   和集合 (日英, 英日)
    if (j か e のどちらかがアラインメント候補で使われていない)
      (j,e) をアラインメント候補に含める
    end
  end
end
```

## 合成の例

(I=積集合, N=隣接点, O=その他の追加, U=削除された点)

	when	the	fluid	pressure	cylinder	31	is	used	,	fluid	is	gradually	applied	.
流体			I											
圧				I										
シリンダ					I									
31						I								
の														
場合	I													
は		N					O							
流体										I				
が											I			
徐々に												I		
排出												U	N	
さ													O	
れる							U						O	
こと													O	
と								O	O					
なる														
。														I

## 合成の例

	this	relation	must	be	maintained	even	after	passing	at	least	100,000	sheets	.
そして							O						
、				O									
上記	I												
関係		I											
を								O					
少なくとも									O	I			
10万											I		
枚											U	N	
通								U			O		
紙												O	
し							U	O					
て						O							
も						I							
維持					I								
し								O					
なけれ			O										
ば			O										
なら			I	U									
ない			O										
。													I

## 合成の例

	first	,	a	description	will	be	given	of	the	structure	of	the	wheels	2	.
まず	I														
、		I													
車輪													I		
2														I	
の								O	N		U				
構造										I					
について					U						N				
説明				O	I	O	O								
する					O							O			
。															I



## 単語単位の翻訳確率

- これまでに，日英，英日の1対 $n$ および $m$ 対1の単語対応を組み合わせて $m$ 対 $n$ の単語対応を得る方法を示した．
- これから，それを利用して，フレーズを抽出する方法を示す．
- その前に，ここで得た $m$ 対 $n$ の対応から単語単位の辞書を作る方法を示す．

## 単語単位の辞書

### アライメント

	e1	e2	e3	e4
j1	*		*	
j2				
j3		*		
j4				*

から ,  $(j1, e1), (j1, e3), (j3, e2), (j4, e4)$  の単語対応がとれるので , これから ,  $j \rightarrow e$  と  $e \rightarrow j$  方向の辞書がつけれる . 更に ,

$$P(j|e) = \frac{\#(j, e)}{\sum_{j'} \#(j', e)}$$

$$P(e|j) = \frac{\#(j, e)}{\sum_{e'} \#(j, e')}$$

により , 単語単位の翻訳確率が求まる . (スムージングをしても良い)

# 単語単位の辞書：頻度 日本語単語 英語単語

2230457 、 the		
1811473 、 ,	146977 さ to	87623 する a
1804260 。 .	146509 2 2	87410 において in
1802874 の of	143482 を for	85936 1 first
1767572 の the	142192 し to	85471 その the
704450 に to	140494 出力 output	84919 として as
668066 は the	139684 制御 control	84888 第 first
584253 が is	135600 が are	84615 に with
543500 は is	130609 は are	84356 メモリ memory
527825 に in	128437 この this	82459 6 6
480009 する the	126068 実施 embodiment	82098 値 value
448118 を is	124007 3 3	81423 に ,
376711 と and	117825 する to	79752 情
376195 、 a	112598 よう as	報 information
365107 図 fig.	111611 た a	79287 1 2 12
345039 を of	110575 ステップ step	79126 1 1 11
287703 ) )	109286 4 4	79048 の for
285186 ( (	106286 および and	78406 例 embodiment
281659 信号 signal	105113 を are	78238 に at
272811 を to	101722 で by	78233 こと can
245721 データ data	101291 及び and	77678 接続 connected
230038 回路 circuit	100440 示す shown	77671 、 an
224630 で in	98899 が of	77513 部 portion
209459 で ,	98684 が can	75575 画像 image
202160 に on	97290 示す in	75000 / /
198721 1 1	96570 5 5	73687 に into
196317 から from	95364 1 0 10	72920 上記 the
193708 この the	95330 電圧 voltage	72590 層 layer
177043 に shown	93903 する be	72466 第 second
175086 し the	93619 入力 input	72399 7 7
166888 、 and	92337 さ in	72231 膜 film
164284 た the	89141 形成 formed	
154881 は a		

# 単語単位の辞書：P(日|英)でソート

emu emu	カテプシン cathepsin	C d x cd.sub.x
S C L K sclk1	W C B R B wcbrb	
+ .sub.+	V p k vpk	x .times..sigma..sub.i3
ウェル -well	V f v vfv	i n o d e inode
L R D lrd	T R B L trbl	S E N P n senpn
L A T lat.sub.	P D R H I T pdrhit	S D B n sdbn
M R S O U T mrsout	C t s c.sub.ts	P F A i pfai
0 0 0 0 0 0 0	4 0 8 0 4080	M O P L S mopls
0 00000000	試験 under-test	K F L E A K M
W D R wdr	r o m a r e a romarea	X kfleakmx
C B S cbs.sub.	R M I S rmis	C A D cad7
P x R pxr	P G pg.sub.	9 8 0 980.degree.
M B T mbt	C T R L R ctrlr	8 5 1 0 8510
M B Q mbq	3 1 1 3.sub.11	8 3 1 83.sub.1
B R E bre	3 0 9 0 3090	1 0 ta.sub.10
6 8 5 685	2 0 4 1 204.sub.1	偏り one-sidedness
. 38q	1 0 0 7 0 10070	微動 micromotion
r a m a r e a ramarea	セル フ プ リ チャー	q se.sub.1-q
k D a kda	ジ self-precharge	m e n t ment
S d w s.sub.dw	w a i t i d waitid	e E C O N O M
R H I T rhit	Z M R zmr	Y eeconomy
C S cs.sub.	V p a s s vpass	T R R Q trrq
d s p l s dspls	O P R H oprh	R S D F rsdf
Z O D L zodl	L I O i lioi	R P R G rprg
G C F gcf	1 9 1 6 1916	Q H N qhn
7 0 3 0 7030	日 date-indicating	I R Q irq9
p r e t r n pretrn	r b l j rblj	H N M O S hnmos
d q m x dqmx	Z C O zco	G L j glj
2 2.sup.	W S E L wsel.sub.	E A V eav
w r e q wreq	V T T A O vttao	D W D E dwde
R W C rwc	R S R a v rsr.sub.av	B j b bjb
S G P sgp	P D R T P i pdrtpi	B R I a bria
3 0 8 0 3080	N s a nsa.sub.	5 0 2 2 5022
排気 exhaling	I a b i.sub.ab	3 1 2 2 3122

# 単語単位の辞書：P(英|日)でソート

オーディオ audio	digest digest	B G M M bgmm
p c h M O S F E	R X T p u l s	8 2 , 8 2 82
T pch-mosfet	e rxtpulse	スレッシュ threshold
I I D R iidr	D R a dra	q c qc
@ @	W B L n wbln	o a m oam
c o a r s e coarse	S P X spx	m e t a p h o
i r o m irom	S E L F O S C selfosc	r metaphor
V c c p vccp	F M T fmt	T H i d l th.sub.idl
G D I gdi	C A S I casi	S P K U spku
M S K i mski	P R E D pred	R M B rmb
C K M ckm	D T G dtg	R F I C rfic
N P R npr	C o p y copy	E X T C L K extclk
S P L spl	C A a caa	B L h blh
R H I T rhit	t P G R tpgr	4 0 9 0 4090
N O X .sigma.nox	d q m x dqmx	n M I S nmis
Z M C H G zmchg	X S U xsu	l u n c h e r luncher
S T X stx	V r c vrc	a u t o m a t o
R I T rit	V c o m H vcomh	n automaton
G B U F R gbufr	S l u slu	V h i g h vhigh
e x t Z R A S extzras	D B B a dbba	S L I F slif
M N P mnp	m p g m mpgm	G W D gwd
コースター coaster	f A fa	G F M gfm
S u p s.sub.up	X B L A xbla	E A D ead
9 1 0 0 9100	R E M E reme	D i n L B n dinlbn
e x t Z C A S extzcas	M B W mbw	C W L cwl
S H L shl	C l t clt	C T R L R ctrlr
I n z in.sub.z	C O C coc	B S B bsb
I m s i.sub.ms	C L E cle	1 , 5 5 5 1,555
B L r blr	R W L E rwle	w a i t i d waitid
フォーマッタ formater	R B L n rbln	v b a t vbat
	N B C nbc	f u p f.sub.up
	B T G btg	X S D xsd
		V Z vz
		S A N G sang
		M I S F E T Q h q.sub.h

# 単語単位の辞書：対数尤度比でソート

。 .		
の of		
、 ’	1 1 11	第 first
、 the	値 value	よう as
の the	この this	電流 current
図 fig.	は the	電極 electrode
信号 signal	層 layer	2 1 21
) )	発明 invention	1 5 15
( (	7 7	本 present
データ data	/ /	動作 operation
に to	画像 image	2 2 22
が is	膜 film	部 portion
回路 circuit	8 8	基板 substrate
と and	方向 direction	半導体 semiconductor
は is	2 0 20	として as
1 1	1 4 14	図 figs.
から from	に on	3 0 30
に in	および and	示す shown
出力 output	する the	トランジスタ transistor
制御 control	1 3 13	例 embodiment
2 2	処理 processing	に shown
3 3	アドレス address	第 second
を is	位置 position	レベル level
4 4	、 a	端子 terminal
電圧 voltage	形成 formed	装置 apparatus
実施 embodiment	及び and	装置 device
1 0 10	で in	を of
ステップ step	ゲート gate	状態 state
5 5	9 9	記録 recording
入力 input	1 6 16	1 first
メモリ memory	接続 connected	ない not
情報 information		
1 2 12		
6 6		

$P(*|\text{流体})$  と  $P(*|\text{fluid})$

# P(\*|流体)

流体 fluid 0.84868421

流体 hydraulic 0.01946272

流体 gas 0.01151316

流体 fluids 0.01096491

流体 flow 0.01041667

流体 liquid 0.00849781

流体 liquid-pressure 0.00603070

流体 hydrodynamic 0.00575658

流体 working 0.00301535

流体 piston 0.00274123

# P(\*|fluid)

fluid 流体 0.45072063

fluid 液 0.23875382

fluid 油 0.10642015

fluid 液体 0.02576794

fluid ブレーキフルード 0.01892561

fluid 流動 0.01368467

fluid 油圧 0.00873490

fluid 流 0.00844373

fluid 作動 0.00698792

fluid 連 0.00684234

## フレーズ対応の抽出

これまでに

	e1	e2	e3	e4
j1	*			
j2		*		
j3		*		
j4			*	*

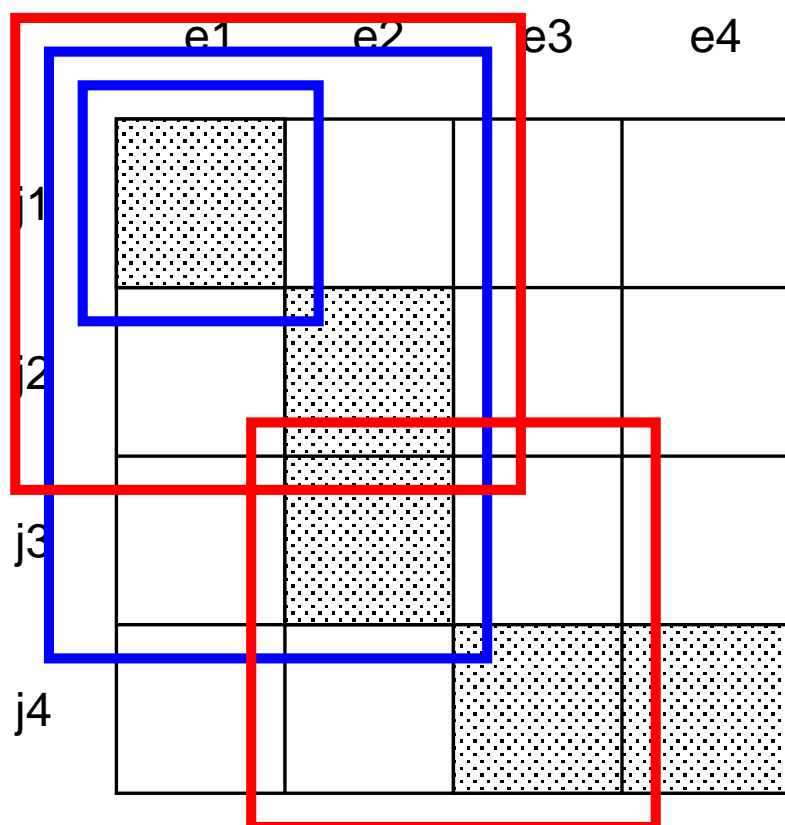
のような単語対応を得る方法を述べた．次は，ここから句対応を得る方法を述べる．

### 基本方針

単語対応に整合するような，なるべく多くの句対応を得る．



## 整合的な句の例



青はOK．赤はダメ．

日本語の句  $J$  と英語の句  $E$  について,  $j \in J$  の対応先を  $e(j)$  とすると,  $e(j) \in E$  でないといけない．

同様に  $e \in E$  の対応先を  $j(e)$  とすると,  $j(e) \in J$  である．  
つまり, 対応の相手は, 日本語句  $J$  内の単語と英語句  $E$  内の単語に限られる．

# 句対応の例：fluid を含む句の例

( 流体 ||| ( fluid ||| 0.166667 0.0307503 1 0.623726 2.718  
 ( 流体 圧 ) ||| ( fluid pressure ) ||| 0.5 0.000742744 1 0.200361 2.718  
 ( 流体 圧 ||| ( fluid pressure ||| 0.5 0.010301 1 0.266855 2.718  
 ( 流体 圧 源 ) ||| ( fluid pressure source ) ||| 1 0.000104349 1 0.159524 2.718  
 ( 流体 圧 源 ) 17 ||| ( fluid pressure source ) 17 ||| 1 9.465e-05 1 0.132599 2.718  
 ( 流体 圧 源 ||| ( fluid pressure source ||| 1 0.00144721 1 0.212466 2.718  
 ) の 送 液 ||| ) of fluid ||| 1 2.88816e-05 1 0.0150938 2.718  
 ) の 送 液 が ||| ) of fluid has ||| 1 6.6955e-06 1 0.000451654 2.718  
 ) の 送 液 が あ っ た ||| ) of fluid has taken ||| 1 2.55555e-09 1 7.43761e-07 2.718  
 、 3 は 液 圧 制 御 装 置 ||| , a fluid pressure control apparatus 3 ||| 1 2.28706e-07 1 0.000292281 2.718  
 、 3 は 加 工 液 ||| , 3 a dielectric fluid ||| 1 3.63074e-05 1 3.96712e-06 2.718  
 、 3 は 加 工 液 、 ||| , 3 a dielectric fluid , and ||| 1 1.47644e-05 1 5.67333e-08 2.718  
 、 3 8 は 加 工 液 ||| , 38 is working fluid ||| 1 0.0088037 1 0.000662696 2.718  
 、 3 8 は 加 工 液 供 給 ||| , 38 is working fluid supplying ||| 1 0.0062226 1 4.34916e-05 2.718  
 、 3 8 は 加 工 液 供 給 手 段 ||| , 38 is working fluid supplying means ||| 1 0.00427184 1 3.16053e-05 2.718  
 、 4 1 d 、 流 体 ||| , 41d , fluid ||| 1 0.0345864 1 0.000467132 2.718  
 、 4 1 d 、 流 体 貯 留 ||| , 41d , fluid storing ||| 1 0.00048108 1 8.14764e-05 2.718  
 、 4 1 d 、 流 体 貯 留 タ ン ク ||| , 41d , fluid storing tank ||| 1 0.000284362 1 7.14906e-05 2.718  
 、 6 お よ び 流 体 ||| and 6 and a fluid ||| 1 0.00410477 1 0.00964129 2.718  
 、 6 お よ び 流 体 測 温 ||| and 6 and a fluid temperature measuring ||| 1 7.86306e-06 1 0.000721089 2.718  
 、 6 1 ' は 伝 熱 流 体 ||| , 61 ' a heat-transfer fluid ||| 1 0.00216443 1 0.000518098 2.718  
 、 6 1 は 伝 熱 流 体 ||| , 61 is a heat-transfer fluid ||| 1 0.0185343 1 0.00119982 2.718  
 、 7 7 は 加 工 液 ||| , 77 is a working fluid ||| 1 0.00864883 1 0.00028376 2.718  
 、 7 7 は 加 工 液 ノ ズ ル ||| , 77 is a working fluid nozzle ||| 1 0.00764329 1 0.000210182 2.718  
 、 7 8 は 加 工 液 ||| , and 78 is a working fluid ||| 1 0.00521192 1 1.09199e-05 2.718  
 、 P は 流 体 ||| , p denotes fluid ||| 1 0.122825 1 0.000660778 2.718  
 、 Q は 流 体 の 流 量 ||| , q denotes flow rate of fluid ||| 1 0.0145985 1 2.77106e-05 2.718  
 、 お よ び 、 第 2 液 ||| and a second fluid ||| 0.333333 9.53798e-05 1 0.00106579 2.718  
 、 お よ び 、 第 2 液 圧 ||| and a second fluid pressure ||| 0.5 3.19513e-05 1 0.00045599 2.718  
 、 この 液 ||| , and this fluid ||| 1 0.0557024 0.166667 0.000360561 2.718  
 、 この 液 ||| , the fluid ||| 0.0285714 0.0047175 0.166667 0.014136 2.718  
 、 この 液 ||| , this fluid ||| 1 0.0934985 0.333333 0.00959799 2.718  
 、 この 液 ||| 8 , and this fluid ||| 1 0.0557024 0.166667 8.11983e-08 2.718  
 、 この 液 圧 ||| , and this fluid pressure ||| 1 0.0186597 0.2 0.000154263 2.718  
 、 この 液 圧 ||| , the fluid pressure ||| 0.25 0.00158031 0.2 0.00604794 2.718  
 、 この 液 圧 ||| , this fluid pressure ||| 1 0.031321 0.4 0.00410641 2.718  
 、 この 液 圧 ||| 8 , and this fluid pressure ||| 1 0.0186597 0.2 3.47399e-08 2.718  
 、 この 液 圧 は ||| , and this fluid pressure is ||| 1 0.00536829 0.333333 4.16064e-05 2.718  
 、 この 液 圧 は ||| , the fluid pressure is ||| 0.5 0.000454647 0.333333 0.0016312 2.718  
 、 この 液 圧 は ||| 8 , and this fluid pressure is ||| 1 0.00536829 0.333333 9.36976e-09 2.718  
 、 この 液 圧 を ||| , this fluid pressure is ||| 1 0.00774038 1 0.000828222 2.718  
 、 この 作 動 液 ||| , the hydraulic fluid ||| 0.125 0.000293068 1 0.000462624 2.718  
 、 この 磁 性 流 体 ||| , the magnetic fluid ||| 0.333333 0.00118769 0.5 0.0798598 2.718  
 、 この 磁 性 流 体 ||| the magnetic fluid ||| 0.0769231 0.000586521 0.5 0.239424 2.718  
 、 この 磁 性 流 体 7 3 ||| , the magnetic fluid 73 ||| 1 0.00108018 1 0.0687977 2.718  
 、 この 流 体 ||| , each fluid ||| 1 0.00229816 1 0.000525922 2.718

# 句対応の例：自動車を含む句の例

- 、 「 自動車 ||| , ‘ automobile ||| 1 0.142122 1 0.0659882 2.718
- 、 「 自動車 電話 ||| , ‘ automobile telephone ||| 1 0.0746708 1 0.0524542 2.718
- 、 「 自動車 電話 」 ||| , ‘ automobile telephone ’ ’ ||| 1 0.0258531 1 0.0386858 2.718
- 、 1 0 2 は 自動車 用 ||| , 102 is a vehicle ||| 1 0.000161113 1 0.00350135 2.718
- 、 1 0 2 は 自動車 用 エンジン ||| , 102 is a vehicle engine ||| 1 9.82163e-05 1 0.0031887 2.718
- 、 3 は 自動車 ||| , 3 a mobile ||| 1 0.000371807 1 0.000168593 2.718
- 、 3 は 自動車 電話 ||| , 3 a mobile telephone ||| 1 0.000195348 1 0.000134015 2.718
- 、 3 は 自動車 電話 局 ||| , 3 a mobile telephone office ||| 1 7.2785e-05 1 3.45847e-06 2.718
- 、 3 は 自動車 電話 局 、 ||| , 3 a mobile telephone office , ||| 1 4.96815e-05 1 1.31658e-06 2.718
- 、 7 は 自動車 ||| , 7 a mobile ||| 1 0.000370597 1 0.000217899 2.718
- 、 7 は 自動車 電話機 ||| , 7 a mobile telephone equipment ||| 1 5.71506e-05 1 4.29458e-07 2.718
- 、 7 は 自動車 電話機 、 ||| , 7 a mobile telephone equipment , ||| 1 3.90098e-05 1 1.63488e-07 2.718
- 、 EC セル 6 を通じて 自動車 ||| vehicle through the ec cell 6 ||| 1 0.00011244 1 0.0167565 2.718
- 、 OA 機器 、 自動車 ||| , business machines , automobiles ||| 1 3.50735e-05 1 2.53658e-07 2.718
- 、 OA 機器 、 自動車 、 ||| , business machines , automobiles and ||| 1 4.58492e-06 1 9.529e-09 2.718
- 、 OA 機器 、 自動車 、 精密 ||| , business machines , automobiles and precision ||| 1 3.50085e-07 1 2.718
- 、 この よう な 電気 自動車 ||| , such an electric vehicle ||| 0.5 1.71207e-05 1 3.34838e-05 2.718
- 、 この よう な 電気 自動車 の ||| , such an electric vehicle ||| 0.5 4.51596e-06 1 3.34838e-05 2.718
- 、 この ハイブリッド 自動車 ||| , in the hybrid vehicle ||| 1 0.000761246 0.5 0.00269411 2.718
- 、 この ハイブリッド 自動車 ||| , in this hybrid vehicle ||| 1 0.0064154 0.5 0.00182924 2.718
- 、 この ハイブリッド 自動車 5 0 0 ||| , in the hybrid vehicle 500 ||| 1 0.00071914 1 0.0022104 2.718
- 、 この ハイブリッド 自動車 5 0 0 において ||| , in the hybrid vehicle 500 , ||| 1 2.0926e-05 1 0.000841 2.718
- 、 この ハイブリッド 自動車 6 0 0 ||| , in this hybrid vehicle 600 ||| 1 0.00591218 1 0.00136645 2.718
- 、 この ハイブリッド 自動車 6 0 0 において ||| , in this hybrid vehicle 600 , ||| 1 0.000172036 1 0.0005 2.718
- 、 しかも 自動車 ||| the automotive ||| 0.1 0.000244648 1 0.0106571 2.718
- 、 たとえば 自動車 用 ||| , for example , an automobile ||| 0.4 0.000664619 1 0.000540882 2.718
- 、 たとえば 自動車 用 エンジン ||| , for example , an automobile engine ||| 1 0.000405159 1 0.000492584 2.718
- 、 コードレス 電話 、 自動車 ||| , cordless telephones , car ||| 1 0.0338496 1 0.000443696 2.718
- 、 コードレス 電話 、 自動車 電話 ||| , cordless telephones , car telephones ||| 1 0.0157655 1 1.50196e-05 2.718
- 、 ディーゼルエンジン 自動車 等 ||| diesel engine automobiles and the like ||| 1 0.0151152 1 3.01743e-05 2.718
- 、 ディーゼルエンジン 自動車 等 に ||| on diesel engine automobiles and the like ||| 1 0.00642812 1 2.24 2.718
- 、 デジタル 自動車 電話 ||| a digital cellular ||| 0.2 2.39509e-05 1 0.000163602 2.718
- 、 デジタル 自動車 電話 システム ||| a digital cellular system ||| 1 9.71561e-06 1 0.000148973 2.718
- 、 デジタル 自動車 ||| the digital cellular ||| 1 5.80871e-05 1 0.00122759 2.718
- 、 デジタル 自動車 電話 ||| the digital cellular telephone ||| 1 3.05191e-05 1 0.000975814 2.718
- 、 デジタル 自動車 電話 5 0 ||| the digital cellular telephone 50 ||| 1 2.83228e-05 1 0.000844553 2.718
- 、 デジタル 自動車 電話 5 0 において ||| the digital cellular telephone 50 , ||| 1 8.24154e-07 1 0.0001 2.718
- 、 ハイブリッド 自動車 ||| a hybrid powered automobile ||| 0.5 0.0354085 0.333333 8.32655e-05 2.718
- 、 ハイブリッド 自動車 ||| the hybrid car ||| 0.166667 0.0441317 0.333333 0.0427639 2.718
- 、 ハイブリッド 自動車 ||| the hybrid vehicle ||| 0.0263158 0.0102417 0.333333 0.0694053 2.718
- 、 ハイブリッド 自動車 5 0 0 ||| the hybrid vehicle 500 ||| 1 0.00967523 1 0.056944 2.718
- 、 ハイブリッド 自動車 5 0 0 の ||| of the hybrid vehicle 500 ||| 1 0.00706622 1 0.0220125 2.718
- 、 ハイブリッド 自動車 5 0 0 の ヨー 方向 ||| yawing directions of the hybrid vehicle 500 ||| 1 0.001016 2.718
- 、 ハイブリッド 自動車 の ||| the hybrid car ||| 0.166667 0.0116407 1 0.0427639 2.718
- 、 ファクシミリ 、 自動車 ||| , facsimiles , automotive vehicles ||| 1 0.0560437 1 1.15431e-06 2.718
- 、 プラスチック 廃棄 物 、 自動車 ||| , plastic wastes , and automobile ||| 1 0.00690582 1 1.86092e-05 2.718

## 現行の句対応抽出の問題点

- 発見的な方法に基づいているため，改良がむずかしい

発見的の意味：抜き出された句対応が良い句対応かどうかの評価が，機械翻訳実験による性能の変化でしか測定できない．確率的なモデルがない．確率的なモデルがあれば，最尤となる句対応を選べば，そのモデルの観点からは最良の句対応が得られる．

- しかし「こういう句対応が良いのでは」といういくつかのモデルは，この発見的な方法と同等か少し性能が落ちる．

## 句の良さの評価

句単位に翻訳していく方法では , なるべく良い句対応を使った翻訳文を得たいので , そのためのスコアを定義する . 一つの句対応について次の5つのスコアが良く使われる .

- 句翻訳確率  $\phi(\bar{f}|\bar{e})$
- 句翻訳確率  $\phi(\bar{e}|\bar{f})$
- 単語確率より  $\text{lex}(\bar{f}|\bar{e})$
- 単語確率より  $\text{lex}(\bar{e}|\bar{f})$
- 句ペナルティ  $\exp(1)$  (固定値)

句翻訳確率  $\phi(\bar{f}|\bar{e})$

$$\phi(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{f}, \bar{e})}{\sum_{\bar{f}'} \text{count}(\bar{f}', \bar{e})}$$

$\text{count}(\bar{f}, \bar{e}) =$  句 $\bar{f}$ と $\bar{e}$ が対応した回数

$\sum_{\bar{f}'} \text{count}(\bar{f}', \bar{e}) =$  句 $\bar{e}$ が出現した回数

$\phi(\bar{f}|\bar{e})$ は，句 $\bar{e}$ が出現したとき，それが $\bar{f}$ に翻訳される確率と考えられる． $\phi(\bar{e}|\bar{f})$ も同様である．

## 単語確率より $\text{lex}(\bar{f}|\bar{e})$

$$\text{lex}(\bar{f}|\bar{e}) = \max_{\mathbf{a}} P_w(\bar{f}|\bar{e}, \mathbf{a})$$

$$P_w(\bar{f}|\bar{e}, \mathbf{a}) = \prod_{i=1}^n E_w(f_i|\bar{e}, \mathbf{a})$$

$$E_w(f_i|\bar{e}, \mathbf{a}) = \frac{1}{|\{j|(i, j) \in \mathbf{a}\}|} \sum_{(i, j) \in \mathbf{a}} w(f_i|e_j)$$

$$w(f_i|e_j) = \text{単語対応の確率} = \frac{\text{count}(f_i, e_j)}{\sum_{f'} \text{count}(f', e_j)}$$

$E_w(f_i|\bar{e}, \mathbf{a}) = w(f_i|e_j)$  の  $e_j$  に関する平均値

$P_w(\bar{f}|\bar{e}, \mathbf{a}) = E_w(f_i|\bar{e}, \mathbf{a})$  の積  $f_i$  の条件付き独立を仮定

$\text{lex}(\bar{f}|\bar{e}) = P_w(\bar{f}|\bar{e}, \mathbf{a})$  が最大の  $\mathbf{a}$  に関する確率を採用

## 句ペナルティ

$\exp(1)$  に固定．スコアとしては，対数を利用するので， $\log(\exp(1)) = 1$  となる．

翻訳に使われた句の  $\log(\text{句ペナルティ})$  の和は，翻訳に使われた句の数である．

実際の翻訳のスコアには，この句ペナルティ(翻訳に使われた句の数)に，ある重み  $\lambda$  を掛けたものを利用する．この  $\lambda$  は自動推定する．

もし， $\lambda > 0$  なら，句の数が多いほど翻訳候補のスコアは良くなる．

$\lambda < 0$  なら，句の数が少ないほど良い．

通常は  $\lambda < 0$  で，少ない方が良い．これは同じ単語数からなる文であれば，使われた句の数が少ない(句の長さは長い)方が良いことになる．



## フレーズ抽出まとめ

発見的手続きを使うことにより，日英のフレーズ対応を得て，そこに各種のスコアを付けることができることをみた．

## 問題 (15分)

自動作成したフレーズテーブルには，さまざまな誤りが含まれる．そのようなフレーズテーブルの中から正しいようなフレーズのみを自動抽出するには，どうしたら良いか．

## 回答例 1

	$\overline{f}$	$\overline{f}$ 以外
$\overline{e}$	a	b
$\overline{e}$ 以外	c	d

$a = \overline{e}$  と  $\overline{f}$  の対応数

$b = \overline{e}$  の出現数  $- a$

$c = \overline{f}$  の出現数  $- a$

$d = \text{句の対応の総数} - a - b - c$

とすると  $\frac{a}{c} \gg \frac{b}{d}$  のとき,  $\overline{e}$  と  $\overline{f}$  は良い対応と考えられる. これは,  $\overline{e}$  と  $\overline{f}$  の独立性の検定をして, 独立性の低いものが良いということなので, その検定結果により, 独立性の低い順に句対応をソートして, 上位のみを利用する.

## 回答例2

- パラレルコーパスを2つに分ける．
- それぞれのコーパスから1つずつフレーズテーブルを作る
- 両方のフレーズテーブルに含まれるようなフレーズのみを採用する

これらの回答例の得失

得 ある程度は信頼性の高いフレーズテーブルができる  
失 カバレッジが低くなる

→ 得が大きいか，失が大きいかは実験してみないと分からない．

今のところ，翻訳の性能を下げずに，ノイズ除去ができるらしいということが分かっている．