

Section 1:

(Question 1A)

- ① Morphology relation of 'talk' and 'talking' is inflection
talking \rightarrow talk + Verb + Progressive
- ② Morphology relation of 'talk' and 'talkative' is derivation
talkative Adj derived from talk Verb
- ③ Morphology relation of 'talk' and ~~talk~~ 'sweet-talk' is:
compound. sweet-talk $>$ sweet + talk

(Section 1B)

~~'the men meet'~~

- ① 'the men meet the men'
~~'the men meet'~~
- ② 'the men meet men'
- ③ 'men meet the men'

(Section 1C) A semantic lexicon defines word senses explicitly.
Here, explicitly defined word senses mean that there are no ambiguity in the word sense like indirect sense, irony, etc.

(Section 1D)

In Lesk algorithm, a 'signature' is a set of words in the gloss and examples of sense. Normally it is from corpus dataset. We can compare signature and the context of the target word, and then to find common words they have. Thus, the 'overlap' can be generated as Overlap(signature, context). If the overlap is greater than the max-overlap we have already known, the best-sense is the corresponding sense.

Xue Yang 19230949

Section 2

(Question 2A)

$$P(\text{sailor} | a) = \frac{c(a \text{ sailor})}{c(a \text{ sailor})} = 1$$

$$P(\text{went} | \text{sailor}) = \frac{c(\text{sailor went})}{c(\text{sailor went})} = 1$$

$$P(\text{to} | \text{went}) = \frac{c(\text{went to})}{c(\text{went to})} = 1$$

$$P(\text{sea} | \text{to}) = \frac{c(\text{to sea})}{c(\text{to sea}) + c(\text{to see})} = 0.5$$

$$P(\text{sea} | \text{sea}) = \frac{c(\text{sea sea})}{c(\text{sea sea}) + c(\text{sea see})} = 1$$

$$P(\text{see} | \text{to}) = \frac{c(\text{to see})}{c(\text{to sea}) + c(\text{to see})} = \frac{1}{2} = 0.5$$

$$P(\text{what} | \text{see}) = \frac{c(\text{see what})}{c(\text{see what}) + c(\text{see see})} = \frac{1}{1+4} = 0.2$$

$$P(\text{he} | \text{what}) = \frac{c(\text{what he})}{c(\text{what he})} = 1$$

$$P(\text{could} | \text{he}) = \frac{c(\text{he could})}{c(\text{he could})} = 1$$

$$P(\text{see} | \text{could}) = \frac{c(\text{could see})}{c(\text{could see})} = 1$$

$$P(\text{see} | \text{see}) = \frac{c(\text{see see})}{c(\text{see see}) + c(\text{see what})} = \frac{4}{4+1} = 0.8$$

$$P(\text{all} | \text{but}) = \frac{c(\text{but all})}{c(\text{but all})} = 1$$

$$P(\text{that} | \text{all}) = 1$$

$$P(\text{he} | \text{that}) = 1$$

$$P(\text{could} | \text{he}) = 1$$

$$P(\text{the} | \text{was}) = \frac{c(\text{was the})}{c(\text{was the})} = 1$$

$$P(\text{bottom} | \text{the}) = \frac{c(\text{the bottom})}{c(\text{the bottom}) + c(\text{the deep})} = \frac{1}{2} = 0.5$$

$$P(\text{of} | \text{bottom}) = 1$$

$$P(\text{the} | \text{of}) = 1$$

$$P(\text{deep} | \text{the}) = \frac{c(\text{the deep})}{c(\text{the deep}) + c(\text{the bottom})} = 0.5$$

Xue Yang 19230949

$$p(\text{blue}|\text{deep}) = 1$$

$$p(\text{sea}|\text{blue}) = 1$$

Thus, we can generate a probability table:

$p(\text{sailor} \text{a}) = 1$	$p(\text{all} \text{but}) = 1$
$p(\text{went} \text{sailor}) = 1$	$p(\text{that} \text{all}) = 1$
$p(\text{to} \text{went}) = 1$	$p(\text{he} \text{that}) = 1$
$p(\text{sea} \text{to}) = 0.5$	$p(\text{could} \text{he}) = 1$
$p(\text{sea} \text{sea}) = 1$	$p(\text{the} \text{was}) = 1$
$p(\text{see} \text{to}) = 0.5$	$p(\text{bottom} \text{the}) = 0.5$
$p(\text{what} \text{see}) = 0.2$	$p(\text{of} \text{bottom}) = 1$
$p(\text{he} \text{what}) = 1$	$p(\text{the} \text{of}) = 1$
$p(\text{could} \text{he}) = 1$	$p(\text{deep} \text{the}) = 0.5$
$p(\text{see} \text{see}) = 0.8$	
$p(\text{see} \text{could}) = 1$	

Given a bigram language model, the formula is:

$$p(w_1, w_2, \dots, w_n) = \prod_{i=1, 2, \dots, n} p(w_i | w_{i-1})$$

And $p(\text{'he could see what'})$

$$= p(\text{he}) \times p(\text{could}|\text{he}) \times p(\text{see}|\text{could}) \times p(\text{what}|\text{see})$$

$$= \frac{2}{33} \times 1 \times 1 \times 0.2 \approx 0.01212$$

Xue Yang 19230949

(Question 2B)

$$p(\text{he could see the sea}) = p(\text{he}) \times p(\text{could}|\text{he}) \times p(\text{see}|\text{could}) \\ \times p(\text{the}|\text{see}) \times p(\text{sea}|\text{the})$$

$$p(\text{he}) = \frac{2+1}{18+33} \quad (\text{the unique words in corpus 23 18}) \\ = \frac{1}{7}$$

$$p(\text{could}|\text{he}) = \frac{c(\text{he could})+1}{c(\text{he could})+18} = \frac{2+1}{2+18} = \frac{3}{20}$$

$$p(\text{see}|\text{could}) = \frac{c(\text{could see})+1}{c(\text{could see})+18} = \frac{2+1}{2+18} = \frac{3}{20}$$

$$p(\text{the}|\text{see}) = \frac{c(\text{see the})+1}{\cancel{c(\text{see the})} + c(\text{see see})+18} = \frac{0+1}{1+4+18} = \frac{1}{23}$$

$c(\text{see what})$

$$p(\text{sea}|\text{the}) = \frac{c(\text{the sea})+1}{c(\text{the deep})+18} = \frac{0+1}{1+18} = \frac{1}{19}$$

$$\text{Thus, } p = \frac{1}{7} \times \frac{3}{20} \times \frac{3}{20} \times \frac{1}{23} \times \frac{1}{19} \approx 0.303 \times 10^{-5}$$

Xue Yang 19230949

(Question 2C)

In machine translation, we wish to get the best translation given a foreign text, say, we wish to find $\text{argmax}(t|f)$ where t is translation and f is foreign text.

Given Bayes's Law:

$$p(t_1|f) = p(f|t_1) \times p(t_1) / p(f)$$

$$p(t_2|f) = p(f|t_2) \times p(t_2) / p(f)$$

It shows that $p(t_1|f) > p(t_2|f)$ if $p(f|t_1)p(t_1) > p(f|t_2)p(t_2)$

Thus, to find the maximum $p(t|f)$ is equivalent to find:

$\text{argmax } p(f|t) * p(t)$, where $p(f|t)$ is translation model, $p(t)$ is language model.

This is how we use language model to resolve ambiguities in machine translation.

Section 3

(Question 3A)

The Hearst patterns are:

① Y including X

② Y such as X

The hyponymy extraction ~~rules~~ with Hearst Patterns rules are:

P5

① IF NP_0 such as $\{NP_1, NP_2, \dots, (\text{and/or})\} NP_n$ THEN for all $NP_i | 1 \leq i \leq n$, hyponymy(NP_i, NP_0)

~~this is a heuristic.~~

① Thus, we can get hyponym { association football, sports }
{ rugby, sports }
(horse racing, sports)
(golf sports)

② 'ZF' 'Y' including $[X_1, \dots, X_n]$ THEN $\{X_1, Y\} \dots \{X_n, Y\}$

Thus, we can get (English language, culture)

(sports, culture)

Xue Yang 19230949

(Question 3B)

Geopolitically (O), Ireland (B-LOC) is (O) divided (O) between (O) the (O) Republic (B-LOC) of (I-LOC) Ireland (Z-LOC) (officially (O) named (O) Ireland (B-LOC)), which (O) covers (O) five-sixths (O) of (O) the (O) island (O), and (O) Northern (B-LOC) Ireland (Z-LOC), which (O) is (O) part (O) of (O) the (O) United (B-LOC) Kingdom (Z-LOC).

The Location entities are: "Ireland", "Republic of Ireland", "Northern Ireland", and "United Kingdom".

According to ZOB tagging scheme, the beginning of the entity is annotated with 'B', the words inside of entity is annotated with 'I', and the words outside of an entity is annotated with 'O'.

(Question 3C)

A knowledge graph can indicate the properties of objects, the relationship between objects, or events. ~~Za~~ The knowledge graphs put data in context via link linking and semantic metadata.

Thus, in addition to textual context information, we can also utilize the links in knowledge graphs to build entity linking. For example, we can use some algorithms like PageRank algorithm to rank pages, or we can generate authority score or hub score (in/out degree) for the link to see how central this entity is.

Section 4

(Question 4A)

"nice" has a positive sentiment score: $1 \times 2 = 2$

"lovely" has a positive sentiment score: $1 \times 2 = 2$

"recommend" has a positive sentiment score: 1

~~Thus, the~~: There are no words that have negative sentiment.

Thus, the Positive (POS) score is 5,

the Negative (NEG) sentiment score is 0.

(Question 4B)

Aspect-based sentiment analysis:

aspect	sentiment
hotel	POS
bar	POS
restaurant	POS
food	POS
staff	POS
hotel	

~~(Question 4c)~~

Question (4C)

~~We need~~ When we do NLP tasks, we need to ~~extract~~ extract information from text and analyse it. For example, we ~~can~~ analyze the data from social media ~~to observe~~ for NLP projects.

But it is important to pay attention to data privacy and data protection issues associated with the data collection for a specific (NLP) task. For example, ~~we~~ sometimes we need to inform the affected individuals what we are doing with their data; we may need to establish the maximum time period required for keeping the data; or we need to develop appropriate data anonymization strategies, etc.

This is because language data is also personal data as individuals can be identified by their language data use. NLP can lead to 'fingerprinting' of individuals. So, data privacy of individuals need to be protected in NLP tasks.