

SIT330-770: Natural Language Processing

Week 5 - Vector Embeddings and Sequence Labeling

Dr. Mohamed Reda Bouadjenek

School of Information Technology, Faculty of Sci Eng & Built Env

reda.bouadjenek@deakin.edu.au

1

SIT330-770: Natural Language Processing

Week 5.1 - Word Meaning

Dr. Mohamed Reda Bouadjenek

School of Information Technology, Faculty of Sci Eng & Built Env

2

What do words mean?

- N-gram or text classification methods we've seen so far
 - Words are just strings (or indices w_i in a vocabulary list)
 - That's not very satisfactory!
- Introductory logic classes:
 - The meaning of "dog" is DOG; cat is CAT
 $\forall x \text{ DOG}(x) \rightarrow \text{MAMMAL}(x)$
- Old linguistic joke by Barbara Partee in 1967:
 - Q: What's the meaning of life?
 - A: LIFE
- That seems hardly better!

3

Desiderata

- What should a theory of word meaning do for us?
- Let's look at some desiderata
- From [lexical semantics](#), the linguistic study of word meaning

4

Lemmas and senses

lemma

mouse (N)

sense

1. any of numerous small rodents...
2. a hand-operated device that controls a cursor...

Modified from the online thesaurus WordNet

A **sense** or "**concept**" is the meaning component of a word
Lemmas can be **polysemous** (have multiple senses)

5

Relations between senses: Synonymy

- Synonyms have the same meaning in some or all contexts.
 - filbert / hazelnut
 - couch / sofa
 - big / large
 - automobile / car
 - vomit / throw up
 - water / H₂O

6

Relations between senses: Synonymy



- Note that there are probably no examples of perfect synonymy.
 - Even if many aspects of meaning are identical
 - Still may differ based on politeness, slang, register, genre, etc.

7

Relation: Synonym?



water/H₂O
 "H₂O" in a surfing guide?
 big/large
 my big sister != my large sister

8

The Linguistic Principle of Contrast



- Difference in form → difference in meaning

9

Abbé Gabriel Girard 1718



Re: "exact" synonyms
 je ne crois pas qu'il y ait de-
 mot synonyme dans aucune
 Langue. Je le dis par con-
 [I do not believe that there
 is a synonymous word in any
 language]



Thanks to Mark Aronoff!

10

Relation: Similarity



Words with similar meanings. Not synonyms, but sharing some
 element of meaning
 car, bicycle
 cow, horse

11

Ask humans how similar 2 words are



word1	word2	similarity
vanish	disappear	9.8
behave	obey	7.3
belief	impression	5.95
muscle	bone	3.65
modest	flexible	0.98
hole	agreement	0.3

SimLex-999 dataset (Hill et al., 2015)

12

Relation: Word relatedness

- Also called "word association"
- Words can be related in any way, perhaps via a semantic frame or field

coffee, tea: **similar**
coffee, cup: **related**, not similar

13

Semantic field

- Words that
 - cover a particular semantic domain
 - bear structured relations with each other.

hospitals
surgeon, scalpel, nurse, anaesthetic, hospital

restaurants
waiter, menu, plate, food, menu, chef

houses
door, roof, kitchen, family, bed

14

Relation: Antonymy

- Senses that are opposites with respect to only one feature of meaning
- Otherwise, they are very similar!

dark/light	short/long	fast/slow	rise/fall
hot/cold	up/down		in/out

- More formally: antonyms can
 - define a binary opposition or be at opposite ends of a scale
 - long/short, fast/slow
- Be reversives:

rise/fall,	up/down
------------	---------

15

Connotation (sentiment)

- Words have **affective** meanings
 - Positive connotations (*happy*)
 - Negative connotations (*sad*)
- Connotations can be subtle:
 - Positive connotation: *copy, replica, reproduction*
 - Negative connotation: *fake, knockoff, forgery*
- Evaluation (sentiment)**
 - Positive evaluation (*great, love*)
 - Negative evaluation (*terrible, hate*)

16

Connotation

Osgood et al. (1957)

- Words seem to vary along 3 affective dimensions:
 - valence**: the pleasantness of the stimulus
 - arousal**: the intensity of emotion provoked by the stimulus
 - dominance**: the degree of control exerted by the stimulus

	Word	Score	Word	Score
Valence	love	1.000	toxic	0.008
	happy	1.000	nightmare	0.005
Arousal	elated	0.960	mellow	0.069
	frenzy	0.965	napping	0.046
Dominance	powerful	0.991	weak	0.045
	leadership	0.983	empty	0.081

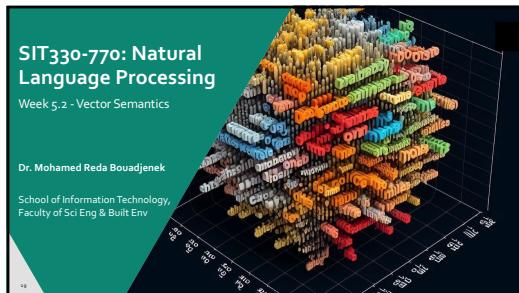
Values from NRC VAD Lexicon (Mohammad 2018)

17

So far

- Concepts** or word senses
 - Have a complex many-to-many association with words (homonymy, multiple senses)
- Have relations with each other
 - Synonymy
 - Antonymy
 - Similarity
 - Relatedness
 - Connotation

18



19

Computational models of word meaning

- Can we build a theory of how to represent word meaning, that accounts for at least some of the desiderata?
- We'll introduce **vector semantics**
 - The standard model in language processing!
 - Handles many of our goals!

20

Ludwig Wittgenstein

• PI #43:
"The meaning of a word is its use in the language"

21

Let's define words by their usages

- One way to define "usage":
◦ words are defined by their environments (the words around them)
- Zellig Harris (1954):
◦ If A and B have almost identical environments we say that they are **synonyms**.

22

What does recent English borrowing *ongchoi* mean?

- Suppose you see these sentences:
 - Ong choi is delicious **sautéed with garlic**.
 - Ong choi is superb **over rice**.
 - Ong choi **leaves** with salty sauces
- And you've also seen these:
 - ...spinach **sautéed with garlic over rice**
 - Chard stems and **leaves** are **delicious**
 - Collard greens and other **salty leafy greens**
- Conclusion:
 - Ongchoi is a leafy green like spinach, chard, or collard greens
 - We could conclude this based on words like "leaves" and "delicious" and "sautéed"

23

Ongchoi: *Ipomoea aquatica* "Water Spinach"

空心菜
kangkong
rau muồng
...

Yamaguchi, Wikimedia Commons, public domain

24

Idea 1: Defining meaning by linguistic distribution

- Let's define the meaning of a word by its distribution in language use, meaning its neighboring words or grammatical environments.

25



Idea 2: Meaning as a point in space (Osgood et al. 1957)

- 3 affective dimensions for a word
 - valence:** pleasantness
 - arousal:** intensity of emotion
 - dominance:** the degree of control exerted

	Word	Score	Word	Score
Valence	love	1.000	toxic	0.008
	happy	1.000	nightmare	0.005
Arousal	elated	0.960	mellow	0.069
	frenzy	0.965	napping	0.046
Dominance	powerful	0.991	weak	0.045
	leadership	0.983	empty	0.081

NRC VAD Lexicon
(Mohammed 2018)

- Hence the connotation of a word is a vector in 3-space

26

26



Idea 1: Defining meaning by linguistic distribution

Idea 2: Meaning as a point in multidimensional space

27

Defining meaning as a point in space based on distribution

- Each word = a vector (not just "good" or " w_{45} ")
- Similar words are "**nearby in semantic space**"
- We build this space automatically by seeing which words are **nearby in text**



28



We define meaning of a word as a vector

- Called an "embedding" because it's embedded into a space (see textbook)
- The standard way to represent meaning in NLP
- Every modern NLP algorithm uses embeddings as the representation of word meaning**
- Fine-grained model of meaning for similarity

29



Intuition: why vectors?

- Consider sentiment analysis:
 - With **words**, a feature is a word identity
 - Feature s_i : "The previous word was 'terrible'"
 - requires **exact same word** to be in training and test
 - With **embeddings**:
 - Feature is a word vector
 - The previous word was vector $[35, 22, 17, \dots]$
 - Now in the test set we might see a similar vector $[34, 21, 14, \dots]$
 - We can generalize to **similar but unseen words!!!**

30



We'll discuss 2 kinds of embeddings

- **tf-idf**
 - Information Retrieval workhorse!
 - A common baseline model
 - Sparse vectors
 - Words are represented by (a simple function of) the **counts** of nearby words
- **Word2vec**
 - Dense vectors
 - Representation is created by training a classifier to **predict** whether a word is likely to appear nearby
 - Later we'll discuss extensions called **contextual embeddings**

31

From now on: Computing with meaning representations instead of string representations

筌者所以在魚，得魚而忘筌 Nets are for fish;
Once you get the fish, you can forget the net.

言者所以在意，得意而忘言 Words are for meaning;
Once you get the meaning, you can forget the words
庄子(Zhuangzi), Chapter 26

32

SIT330-770: Natural Language Processing
Week 5.3 – Words and Vectors: BOW

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env



33

Bag of Words

- A document is represented as vector of words.
- One dimension per word.
- Vector size is the vocabulary size, e.g., English may contain 100k words.
- Different weighting schemas can be used, e.g., tf, log(tf), tf-idf, Boolean, etc.
- Sparse vector, e.g., almost all values are zeros.

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	14	80	62	89
fool	36	58	1	4
wit	20	15	2	3

34

Order matters for NLP tasks!

- Assumes independence between words:
 - The sentences "**John likes Mary**" has the same representation as "**Mary likes John**" – even though the semantic is different).
- May work well for Information Retrieval tasks, but not for NLP tasks!
 - Sentiment analysis:
"Ah **no**, there are good movies on Netflix!" vs. "Ah, there are **no** good movies on Netflix!"

35

Computing word similarity: Dot product and cosine

- The dot product between two vectors is a scalar:
$$\text{dot product}(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^N v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$$
- The dot product tends to be high when the two vectors have large values in the same dimensions
- Dot product can thus be a useful similarity metric between vectors

36

Problem with raw dot-product

- Dot product favors long vectors
- Dot product is higher if a vector is longer (has higher values in many dimension)
- Vector length:

$$\|\mathbf{v}\| = \sqrt{\sum_{i=1}^N v_i^2}$$

- Frequent words (of, the, you) have long vectors (since they occur many times with other words).
- So dot product overly favors frequent words

37

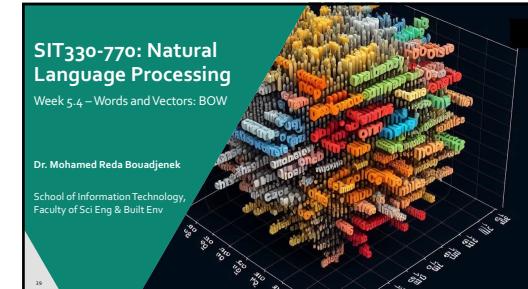
Alternative: cosine for computing word similarity

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

Based on the definition of the dot product between two vectors \mathbf{a} and \mathbf{b}

$$\frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} = \frac{|\mathbf{a}| |\mathbf{b}| \cos \theta}{|\mathbf{a}| |\mathbf{b}|} = \cos \theta$$

38



Bag of Words

- A document is represented as vector of words.
- One dimension per word.
- Vector size is the vocabulary size, e.g., English may contain 100k words.
- Different weighting schemas can be used, e.g., tf, log(tf), tf-idf, Boolean, etc.
- Sparse vector, e.g., almost all values are zeros.

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	14	80	62	89
fool	36	58	1	4
wit	20	15	2	3

40

Order matters for NLP tasks!

- Assumes independence between words:
 - The sentences "**John likes Mary**" has the same representation as "**Mary likes John**" – even though the semantic is different.
- May work well for Information Retrieval tasks, but not for NLP tasks!
 - Sentiment analysis:
"Ah **no**, there are good movies on Netflix!" vs. "Ah, there are **no** good movies on Netflix!"

41

Computing word similarity: Dot product and cosine

- The dot product between two vectors is a scalar:
$$\text{dot product}(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^N v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$$
- The dot product tends to be high when the two vectors have large values in the same dimensions
- Dot product can thus be a useful similarity metric between vectors

42

Problem with raw dot-product

- Dot product favors long vectors
- Dot product is higher if a vector is longer (has higher values in many dimension)
- Vector length:

$$\|\mathbf{v}\| = \sqrt{\sum_{i=1}^N v_i^2}$$
- Frequent words (of, the, you) have long vectors (since they occur many times with other words).
- So dot product overly favors frequent words

43

Alternative: cosine for computing word similarity

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

Based on the definition of the dot product between two vectors a and b

$$\frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} = \frac{|\mathbf{a}| |\mathbf{b}| \cos \theta}{|\mathbf{a}| |\mathbf{b}|} = \cos \theta$$

44

Sparse versus dense vectors

- tf-idf (or PMI) vectors are
 - **long** (length $|V|= 20,000$ to $50,000$)
 - **sparse** (most elements are zero)
- Alternative: learn vectors which are
 - **short** (length $50\text{-}1000$)
 - **dense** (most elements are non-zero)

45

Sparse versus dense vectors

- Why dense vectors?
 - Short vectors may be easier to use as **features** in machine learning (fewer weights to tune)
 - Dense vectors may **generalize** better than explicit counts
 - Dense vectors may do better at capturing synonymy:
`car` and `automobile` are synonyms; but are distinct dimensions
 - a word with `car` as a neighbor and a word with `automobile` as a neighbor should be similar, but aren't
 - In practice, they work better

46

Common methods for getting short dense vectors

- "Neural Language Model"-inspired models
 - Word2vec (skipgram, CBOW), GloVe
- Singular Value Decomposition (SVD)
 - A special case of this is called LSA – Latent Semantic Analysis
- Alternative to these "static embeddings":
 - Contextual Embeddings (ELMo, BERT)
 - Compute distinct embeddings for a word in its context
 - Separate embeddings for each token of a word

47

Simple static embeddings you can download!

- Word2vec (Mikolov et al)
- <https://code.google.com/archive/p/word2vec/>
- GloVe (Pennington, Socher, Manning)
 - <http://nlp.stanford.edu/projects/glove/>

48

Word2vec

- Popular embedding method
- Very fast to train
- Code available on the web
- Idea: predict rather than count
- Word2vec provides various options. We'll do:
 - skip-gram with negative sampling (SGNS)

49

Word2Vec

- Instead of counting how often each word w occurs near "apricot"
 - Train a classifier on a binary prediction task:
 - Is w likely to show up near "apricot"?
- We don't actually care about this task
 - But we'll take the learned classifier weights as the word embeddings
- Big idea: self-supervision:
 - A word c that occurs near apricot in the corpus acts as the gold "correct answer" for supervised learning
 - No need for human labels
 - Bengio et al. (2003); Collobert et al. (2011)

50

Approach: predict if candidate word c is a "neighbor"

- Treat the target word t and a neighboring context word c as positive examples.
- Randomly sample other words in the lexicon to get negative examples
- Use logistic regression to train a classifier to distinguish those two cases
- Use the learned weights as the embeddings

51

Skip-Gram Training Data

- (assuming a +/- 2 word window)
...lemon, a [tablespoon of **apricot** jam, a] pinch...
- $c_1 \quad c_2 \quad [target] \quad c_3 \quad c_4$
- Goal: train a classifier that is given a candidate (word, context) pair
 - (apricot, jam)
 - (apricot, aardvark)
 - ...
- And assigns each pair a probability:
 - $P(+|w, c)$
 - $P(-|w, c) = 1 - P(+|w, c)$

52

Similarity is computed from dot product

- Remember: two vectors are similar if they have a high dot product
 - Cosine is just a normalized dot product
- So:
 - Similarity(w, c) $\propto w \cdot c$
- We'll need to normalize to get a probability
 - (cosine isn't a probability either)

53

Turning dot products into probabilities

- $\text{Sim}(w, c) = w \cdot c$
- To turn this into a probability
- We'll use the sigmoid from logistic regression:

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

$$P(-|w, c) = 1 - P(+|w, c)$$

$$= \sigma(-c \cdot w) = \frac{1}{1 + \exp(c \cdot w)}$$

54

How Skip-Gram Classifier computes $P(+|w, c)$

$$P(+|w, c) = \sigma(c \cdot w) = \frac{1}{1 + \exp(-c \cdot w)}$$

- This is for one context word, but we have lots of context words.
- We'll assume independence and just multiply them:

$$P(+|w, c_{1:L}) = \prod_{i=1}^L \sigma(c_i \cdot w)$$

$$\log P(+|w, c_{1:L}) = \sum_{i=1}^L \log \sigma(c_i \cdot w)$$

55

Skip-gram classifier: summary

- A probabilistic classifier, given
 - a test target word w
 - its context window of L words $c_{1:L}$
- Estimates probability that w occurs in this window based on similarity of w (embeddings) to $c_{1:L}$ (embeddings).
- To compute this, we just need embeddings for all the words.

56

SIT330-770: Natural Language Processing

Week 5.5 – Word2vec: Learning the embeddings

Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

57

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

C1	C2 [target]	C3	C4
t	c		

positive examples +

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

58

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

C1	C2 [target]	C3	C4
t	c		

positive examples +

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

For each positive example we'll grab k negative examples, sampling by frequency

59

Skip-Gram Training data

...lemon, a [tablespoon of apricot jam, a] pinch...

C1	C2 [target]	C3	C4
t	c	t	c

positive examples +

t	c
apricot	tablespoon
apricot	of
apricot	jam
apricot	a

negative examples -

t	c
apricot	aardvark
apricot	my
apricot	where
apricot	coaxial
apricot	seven
apricot	dear
apricot	if

60

Word2vec: how to learn vectors

- Given the set of positive and negative training instances, and an initial set of embedding vectors
- The goal of learning is to adjust those word vectors such that we:
 - Maximize the similarity of the target word, context word pairs (w, c_{pos}) drawn from the positive data
 - Minimize the similarity of the (w, c_{neg}) pairs drawn from the negative data.

61

Loss function for one w with $c_{pos}, c_{neg_1}, \dots, c_{neg_k}$

- Maximize the similarity of the target with the actual context words, and minimize the similarity of the target with the k negative sampled non-neighbor words.

$$\begin{aligned} L_{CE} &= -\log \left[P(+|w, c_{pos}) \prod_{i=1}^k P(-|w, c_{neg_i}) \right] \\ &= -\left[\log P(+|w, c_{pos}) + \sum_{i=1}^k \log P(-|w, c_{neg_i}) \right] \\ &= -\left[\log P(+|w, c_{pos}) + \sum_{i=1}^k \log (1 - P(+|w, c_{neg_i})) \right] \\ &= -\left[\log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right] \end{aligned}$$

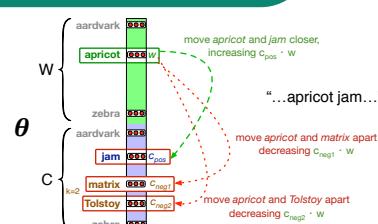
62

Learning the classifier

- How to learn?
 - Stochastic gradient descent!
- We'll adjust the word weights to
 - make the positive pairs more likely
 - and the negative pairs less likely,
 - over the entire training set.

63

Intuition of one step of gradient descent



64

Reminder: gradient descent

- At each step
 - Direction: We move in the reverse direction from the gradient of the loss function
 - Magnitude: we move the value of this gradient $\frac{d}{dw} L(f(x; w), y)$ weighted by a learning rate η
 - Higher learning rate means move w faster

$$w^{t+1} = w^t - \eta \frac{d}{dw} L(f(x; w), y)$$

65

The derivatives of the loss function

$$\begin{aligned} L_{CE} &= -\left[\log \sigma(c_{pos} \cdot w) + \sum_{i=1}^k \log \sigma(-c_{neg_i} \cdot w) \right] \\ \frac{\partial L_{CE}}{\partial c_{pos}} &= [\sigma(c_{pos} \cdot w) - 1]w \\ \frac{\partial L_{CE}}{\partial c_{neg_i}} &= [\sigma(c_{neg_i} \cdot w)]w \\ \frac{\partial L_{CE}}{\partial w} &= [\sigma(c_{pos} \cdot w) - 1]c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i} \cdot w)]c_{neg_i} \end{aligned}$$

66

Update equation in SGD

- Start with randomly initialized C and W matrices, then incrementally do updates

$$c_{pos}^{t+1} = c_{pos}^t - \eta [\sigma(c_{pos}^t \cdot w^t) - 1] w^t$$

$$c_{neg}^{t+1} = c_{neg}^t - \eta [\sigma(c_{neg}^t \cdot w^t)] w^t$$

$$w^{t+1} = w^t - \eta \left[[\sigma(c_{pos} \cdot w^t) - 1] c_{pos} + \sum_{i=1}^k [\sigma(c_{neg_i} \cdot w^t)] c_{neg_i} \right]$$

67

Two sets of embeddings

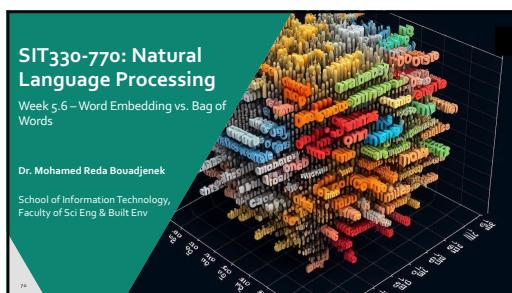
- SGNS learns two sets of embeddings
 - Target embeddings matrix W
 - Context embedding matrix C
- It's common to just add them together, representing word i as the vector $w_i + c_i$

68

Summary: How to learn word2vec (skip-gram) embeddings

- Start with V random d-dimensional vectors as initial embeddings
- Train a classifier based on embedding similarity
 - Take a corpus and take pairs of words that co-occur as positive examples
 - Take pairs of words that don't co-occur as negative examples
 - Train the classifier to distinguish these by slowly adjusting all the embeddings to improve the classifier performance
 - Throw away the classifier code and keep the embeddings.

69



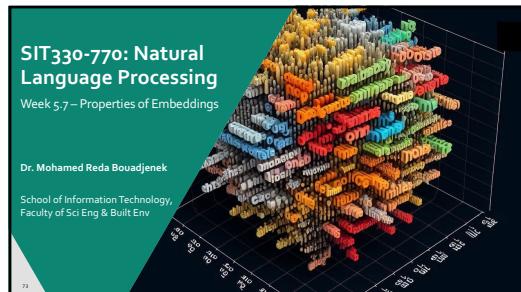
70

Traditional Method - Bag of Words Model	Word Embeddings
Two approaches:	
<ul style="list-style-type: none"> Either sparse one-hot encoding. <ul style="list-style-type: none"> Each word in the vocabulary is represented by one bit position in a HUGE vector. For example, if we have a vocabulary of 10,000 words, and "aaardvark" is the 4th word in the dictionary, it would be represented by [0,0,0,1,0, ..., 0,0]. Or using document representation. <ul style="list-style-type: none"> Each word in the vocabulary is represented by its presence in documents. For example, if we have a corpus of 1M documents, and "Hello" appears in 100,000 documents only, it would be represented by [1,0,1,0, ..., 0,0]. Assumes independence between words. 	<ul style="list-style-type: none"> Stores each word in as a point in space, where it is represented by a dense vector of fixed number of dimensions (generally 300). Dimensions are projections along different axes, more of a mathematical concept. Unsupervised, built just by reading huge corpus. Assumes dependence between words.

71

Traditional Method - Bag of Words Model	Word Embeddings
<ul style="list-style-type: none"> Requires very large weight matrix for 1st layers. <ul style="list-style-type: none"> 10,000 words → 100 units W's size is $10,000 \times 100 = 10^6$ Models are flexible with unseen words in the training set. <ul style="list-style-type: none"> LM: He → is → a → cultivator (X) LM: He → is → a → farmer (≈) 	<ul style="list-style-type: none"> A compact weight matrix for 1st layers. <ul style="list-style-type: none"> d300 → 100 units W's size is $300 \times 100 = 3 \times 10^4$ Flexible models with unseen words in the training set. <ul style="list-style-type: none"> LM: He → is → a → cultivator (X) LM: He → is → a → farmer (≈)

72



73

The kinds of neighbors depend on window size

- Small windows ($C = +/- 2$)**: nearest words are syntactically similar words in same taxonomy
 - Hogwarts* nearest neighbors are other fictional schools
 - Sunnydale, Evernight, Blandings*
- Large windows ($C = +/- 5$)**: nearest words are related words in same semantic field
 - Hogwarts* nearest neighbors are Harry Potter world:
 - Dumbledore, half-blood, Malfoy*

74

Analogical relations

- The classic parallelogram model of analogical reasoning (Rumelhart and Abrahamson 1973)
- To solve: "*apple* is to *tree* as *grape* is to _____"
- Add *tree - apple* to *grape* to get *vine*

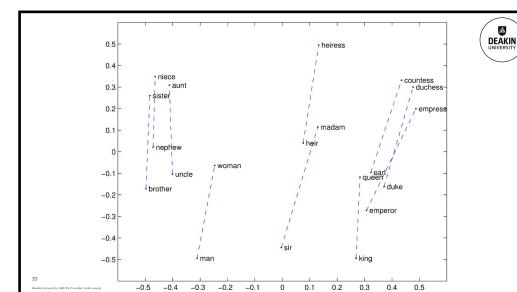
75

Analogical relations via parallelogram

- The parallelogram method can solve analogies with both sparse and dense embeddings (Turney and Littman 2005, Mikolov et al. 2013b)
- king – man + woman is close to queen
- Paris – France + Italy is close to Rome
- For a problem $a:a^*:b:b^*$, the parallelogram method is:

$$\hat{b}^* = \operatorname{argmax}_x \text{distance}(x, a^* - a + b)$$

76



77

Caveats with the parallelogram method

- It only seems to work for frequent words, small distances and certain relations (relating countries to capitals, or parts of speech), but not others. (Linzen 2016, Gladkova et al. 2016, Ethayarajh et al. 2019a)
- Understanding analogy is an open area of research (Peterson et al. 2020)

78

Embeddings as a window onto historical semantics

- Train embeddings on different decades of historical text to see meanings shift
~30 million books, 1850-1990, Google Books data

William L. Hamilton, Jure Leskovec, and Dan Jurafsky. 2016. Diachronic Word Embeddings Reveal Statistical Laws of Semantic Change. Proceedings of ACL.

79

Embeddings reflect cultural bias!

- Ask "Paris : France :: Tokyo : x"
x = Japan
- Ask "father : doctor :: mother : x"
x = nurse
- Ask "man : computer programmer :: woman : x"
x = homemaker

Algorithms that use embeddings as part of e.g., hiring searches for programmers, might lead to bias in hiring

80

Historical embedding as a tool to study cultural biases

Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? Debasing word embeddings." In NeurIPS, pp. 4349-4357. 2016.

- Compute a **gender or ethnic bias** for each adjective: e.g., how much closer the adjective is to "woman" synonyms than "man" synonyms, or names of particular ethnicities
- Embeddings for **competence** adjective (*smart, wise, brilliant, resourceful, thoughtful, logical*) are biased toward men, a bias slowly decreasing 1960-1990
- Embeddings for **dehumanizing** adjectives (*barbaric, monstrous, bizarre*) were biased toward Asians in the 1930s, bias decreasing over the 20th century.
- These match the results of old surveys done in the 1930s

81

SIT330-770: Natural Language Processing

Week 5 - Sequence Labeling

Dr. Mohamed Reda Bouadjenek

School of Information Technology, Faculty of Sci Eng & Built Env

reda.bouadjenek@deakin.edu.au

82

SIT330-770: Natural Language Processing

Week 5.8 - English Word Classes

Dr. Mohamed Reda Bouadjenek

School of Information Technology, Faculty of Sci Eng & Built Env

83

Parts of Speech

- From the earliest linguistic traditions (Yaska and Panini 5th C. BCE, Aristotle 4th C. BCE), the idea that words can be classified into grammatical categories
 - part of speech, word classes, POS, POS tags
- 8 parts of speech attributed to Dionysius Thrax of Alexandria (c. 1st C. BCE):
 - noun, verb, pronoun, preposition, adverb, conjunction, participle, article
- These categories are relevant for NLP today.

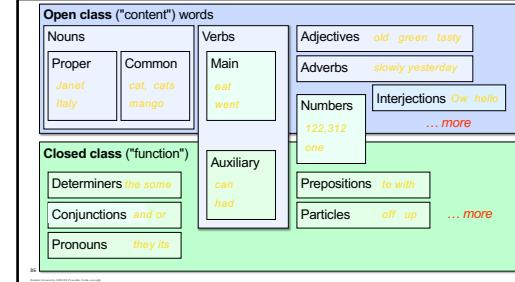
84

Two classes of words: Open vs. Closed

- Closed class words**
 - Relatively fixed membership
 - Usually function words: short, frequent words with grammatical function
 - determiners: *a, an, the*
 - pronouns: *she, he, I*
 - prepositions: *on, under, over, near, by, ...*
- Open class words**
 - Usually content words: Nouns, Verbs, Adjectives, Adverbs
 - Plus interjections: *oh, ouch, uh-huh, yes, hello*
 - New nouns and verbs like *iPhone* or *to fax*



85



Part-of-Speech Tagging

- Assigning a part-of-speech to each word in a text.
- Words often have more than one POS.
- book:**
 - VERB: (*Book that flight*)
 - NOUN: (*Hand me that book*).



87

"Universal Dependencies" Tagset

Tag	Description	Example
Open Class		
ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
VERB	words for actions, events, states, etc.	<i>algorithm, cat, mango, beauty draw, practice, go</i>
PROPN	Proper name: name of a person, organization, place, etc.	<i>Russia, IBM, Colorado oh, um, yet, before</i>
INTJ	Interjection: exclamation, greeting, yes/no response, etc.	
Closed Class Words		
ADP	Adposition (Preposition/Postposition): marks a noun's spatial, temporal, or other relation	<i>in, on, by under</i>
CCONJ	Coordinating Conjunction: joins two phrases/classes	<i>can, nor, should, are and, or, but, either</i>
DET	Determiner: marks noun phrase properties	<i>a, an, the, this one, two, first, second</i>
NUM	Numerical	<i>up, down, on, off, in, out, at, by</i>
PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others that, which</i>
SCONJ	Subordinating Conjunction: joins a main clause with a relative clause such as a sentential complement	
Other		
PUNCT	Punctuation	<i>: ; , . ! ?</i>
SYM	Symbols like \$ or emoji	<i>\$, % asdf, qwfg</i>
X	Other	

Nivre et al. 2016



88

Sample "Tagged" English sentences

- There/**PRO** were/**VERB** 70/**NUM** children/**NOUN** there/**ADV** ./**PUNC**
- Preliminary/**ADJ** findings/**NOUN** were/**AUX** reported/**VERB** in/**ADP** today/**NOUN**'s/**PART** New/**PROPN** England/**PROPN** Journal/**PROPN** of/**ADP** Medicine/**PROPN**



89

SIT330-770: Natural Language Processing

Week 5.9 - Part of Speech Tagging

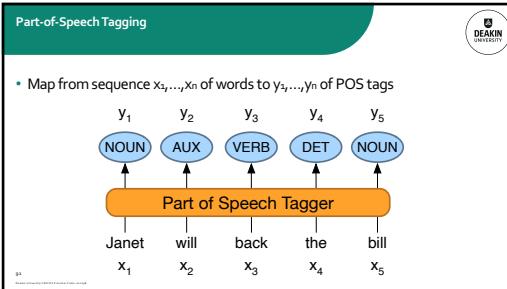
Dr. Mohamed Reda Bouadjenek

School of Information Technology,
Faculty of Sci Eng & Built Env

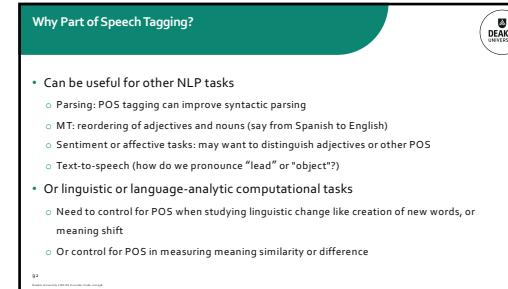


LanguageTool is your intelligent checker for all common errors and language variants in your text. It highlights misspellings, grammar errors, punctuation mistakes and other problems. Errors will not be underlined in the original text, but will be highlighted in red. Grammar errors are highlighted in green. If you click on a red error message, it will be explained in a detailed manner by underlining them in blue. If you like to check a word, just type it in the search bar and click "Search".

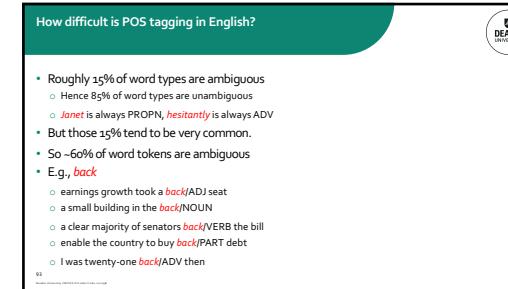
80



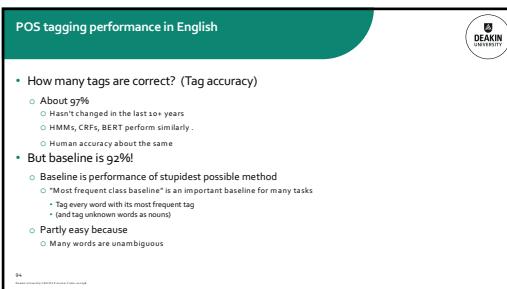
91



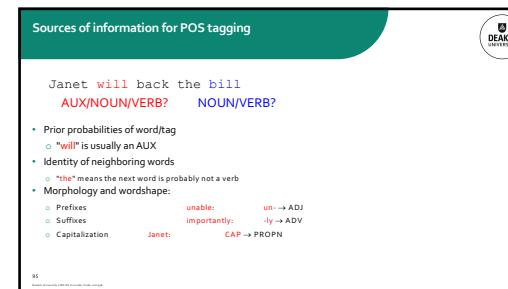
92



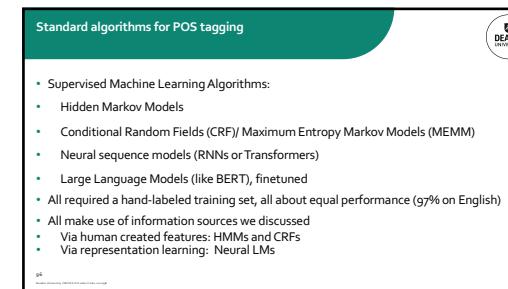
93



94



95



96

SIT330-770: Natural Language Processing
Week 5.10 - Named Entity Recognition (NER)
Dr. Mohamed Reda Bouadjenek
School of Information Technology,
Faculty of Sci Eng & Built Env

97

Named Entities

- Named entity, in its core usage, means anything that can be referred to with a proper name. Most common 4 tags:
 - PER (Person): "Marie Curie"
 - LOC (Location): "NewYorkCity"
 - ORG (Organization): "StanfordUniversity"
 - GPE (Geo-Political Entity): "Boulder, Colorado"
- Often multi-word phrases
- But the term is also extended to things that aren't entities:
 - dates, times, prices

98

Named Entity tagging

- The task of named entity recognition (NER):
 - find spans of text that constitute proper names
 - tag the type of the entity.

99

NER output

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

100

Why NER?

- Sentiment analysis: consumer's sentiment toward a particular company or person?
- Question Answering: answer questions about an entity?
- Information Extraction: Extracting facts about entities from text.

101

Why NER is hard

- Segmentation
 - In POS tagging, no segmentation problem since each word gets one tag.
 - In NER we have to find and segment the entities!
- Type ambiguity

[PER Washington] was born into slavery on the farm of James Burroughs.
[ORG Washington] went up 2 games to 1 in the four-game series.
Blair arrived in [LOC Washington] for what may well be his last state visit.
In June, [GPE Washington] passed a primary seatbelt law.

102

BIO Tagging

- How can we turn this structured problem into a sequence problem like POS tagging, with one label per word?
- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago] route.

Now we have one tag per token!!!

103

BIO Tagging

- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago] route.

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
-	O

Now we have one tag per token!!!

104

BIO Tagging

B: token that *begins* a span
I: tokens *inside* a span
O: tokens outside of any span

of tags (where n is #entity types):
 $1 \text{ O tag},$
 $n \text{ B tags},$
 $n \text{ I tags}$
 $\text{total of } 2n+1$

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
-	O

105

BIO Tagging variants: IO and BIOES

- [PER Jane Villanueva] of [ORG United], a unit of [ORG United Airlines Holding], said the fare applies to the [LOC Chicago] route.

Words	IO Label	BIO Label	BIOES Label
Jane	I-PER	B-PER	B-PER
Villanueva	I-PER	I-PER	E-PER
of	O	O	O
United	I-ORG	B-ORG	B-ORG
Airlines	I-ORG	I-ORG	I-ORG
Holding	I-ORG	I-ORG	E-ORG
discussed	O	O	O
the	O	O	O
Chicago	I-LOC	B-LOC	S-LOC
route	O	O	O
-	O	O	O

106

Standard algorithms for NER

- Supervised Machine Learning given a human-labeled training set of text annotated with tags
 - Hidden Markov Models
 - Conditional Random Fields (CRF)/Maximum Entropy Markov Models (MEMM)
 - Neural sequence models (RNNs or Transformers)
 - Large Language Models (like BERT), finetuned

107

SIT330-770: Natural Language Processing

Week 5.11 – Hidden Markov Model (HMM) Part-of-Speech Tagging

Dr. Mohamed Reda Bouadjenek
School of Information Technology, Faculty of Sci Eng & Built Env

LanguageTool is your intelligent grammar checker for all common browsers and your text box. It highlights misspelled words and grammar errors. Errors will be underlined in red, suggesting corrections with an alternative. Punctuation grammar errors are highlighted in green. You can report errors directly to the LanguageTool community. Click on a suggestion to see its details. You can further elaborate on the error by clicking on the 'Report Error' link. If you like to use LanguageTool in a mobile manner by underscoring them in blue, do you know how? Just click on a word! It's impressively versatile tool, e.g. if you like to know about what happened at 5 PM in the afternoon on Monday, 21 May 2007, by double-clicking a word!

108

Introduction to Markov Chains

- A Markov chain models the probabilities of state sequences, each drawn from a specific set.
- It assumes the future state depends only on the current state, not any prior ones.
- Markov chains are used to predict various phenomena
 - E.g., modeling weather patterns or word sequences.

Figure 8.1 A Markov chain for weather

109

Markov Chain Representation

$Q = q_1, q_2, \dots, q_N$ a set of N states
 $A = a_{11}, a_{12}, \dots, a_{N1}, \dots, a_{NN}$ a transition probability matrix A , each a_{ij} represents the probability of moving from state i to state j , s.t.
 $\sum_j a_{ij} = 1, \forall i$
 $\pi = \pi_1, \pi_2, \dots, \pi_N$ an initial probability distribution over states. π_i is the probability that the Markov chain will start in state i . Some states j may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^N \pi_i = 1$

Figure 8.2 A Markov chain for weather (a) and one for words (b). Observed states and transitions are shown in blue, while hidden states and transitions are shown in grey. (a) shows a 4-state Markov chain for weather, starting at state 1 (SUNNY) with probability 0.5 of starting in state 2 (CLOUDY), probability 0.1 of starting in state 3 (RAINY), etc.

- Markov Assumption:**
 - Formally stated as: $P(q_i=a|q_1, \dots, q_{i-1}) = P(q_i=a|q_{i-1})$ implying that when predicting the future, only the present state matters

110

The Hidden Markov Model

- A Markov chain computes probabilities for sequences of observable events.
- But often, the events of interest are hidden.
 - Example:** Part-of-speech tags in text—hidden because we don't observe them directly.
- Solution:** Hidden Markov Model (HMM) handles both observed and hidden events.
 - HMMs augment Markov chains

111

Probabilistic Sequence Modeling with HMMs

- A Hidden Markov Models (HMM) is a probabilistic sequence model that, given a sequence of units (words, letters, morphemes, sentences, etc.), computes a probability distribution over possible sequences of labels.
 - HMMs determine the likelihood of different label sequences and select the most probable sequence based on the observed data.
 - HMM is based on augmenting the Markov chain

112

Input and Assumptions

- Input (O):** Sequence of observations (o_1, o_2, \dots, o_T) drawn from vocabulary V .
- Assumptions of first-order HMM:**
 - Markov Assumption:**
 - Probability of state q_t depends only on the previous state (q_{t-1}) .
 - $P(q_t|q_{t-1}, \dots, q_1) = P(q_t|q_{t-1})$
 - Output Independence:**
 - Probability of observation o_t depends only on the state that produced it q_t .
 - $P(o_t|q_{t-1}, \dots, q_1, o_1, \dots, o_{t-1}, o_t) = P(o_t|q_t)$

113

SIT330-770: Natural Language Processing

Week 5.12 – The components of an HMM tagger

Dr. Mohamed Reda Bouadjenek
 School of Information Technology,
 Faculty of Sci Eng & Built Env

114

Components of an HMM Tagger

- An HMM tagger consists of two main components:
 - Matrix A which represents the tag transition.
 - Matrix B which represents emission probabilities.

115

The A Matrix - Transition Probabilities

- The A matrix encapsulates the tag transition probabilities, $P(t|t_{-1})$, which express how likely a tag follows its predecessor.
- Example:
 - The modal verb "will" commonly precedes the base form of a verb (VB), as in "will race", leading to a high transition probability.
 - These probabilities are derived using maximum MLE by counting tag occurrences in a labeled corpus.
- Calculating Transition Probabilities:
 - In the WSJ corpus example, the modal verb tag (MD) is observed 13,124 times.
 - Out of these, MD transitions to a base verb (VB) 10,471 times.
 - Using MLE, we estimate $P(VB|MD) = C(MD, VB) / C(MD) = 10,471 / 13,124 = 0.31$.

116

The B Matrix - Emission Probabilities

- The B matrix contains emission probabilities, $P(w|t)$, which quantify the likelihood of a word being tagged with a specific tag.
- Emission Probability Calculation
 - To calculate emission probabilities, we count how often a word occurs with a particular tag in a corpus.
 - For instance, the MD tag associated with the word 'will' occurs 4,046 times in the WSJ corpus.
 - Hence, $P(will|MD)$ is calculated as $C(MD, will) / C(MD) = 4,046 / 13,124 = 0.31$.

117

Components of HMM

$Q = q_1 q_2 \dots q_N$ a set of N states
 $A = a_{ij} \dots a_{ij} \dots a_{NN}$ a **transition probability matrix** A , each a_{ij} representing the probability of moving from state i to state j , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$
 $B = b_i(o_t)$
 $\pi = \pi_1, \pi_2, \dots, \pi_N$ a sequence of observation likelihoods, also called **emission probabilities**, each expressing the probability of an observation o_t (drawn from a vocabulary $V = v_1, v_2, \dots, v_V$) being generated from a state q_i , an **initial probability distribution** over states. π_i is the probability that the Markov chain will start in state i . Some states j may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^N \pi_i = 1$

118

SIT330-770: Natural Language Processing

Week 5.13 – HMM tagging as decoding

Dr. Mohamed Reda Bouadjenek

School of Information Technology,
Faculty of Sci Eng & Built Env

119

Decoding with Hidden Markov Models

- Decoding is the process of determining the most probable sequence of hidden states (tags) based on observed data.
 - Given a sequence of observations $O = o_1, o_2, \dots, o_r$, decoding aims to find the most probable sequence of states $Q = q_1, q_2, \dots, q_r$.
 - The input is an HMM $\lambda = (A, B)$, with A being the transition probabilities and B the emission probabilities.

$$\hat{t}_{1:n} = \underset{t_1 \dots t_n}{\operatorname{argmax}} P(t_1 \dots t_n | w_1 \dots w_n)$$

120

Decoding with Hidden Markov Models (i)

$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n)$ MAP is "maximum a posteriori" = most likely sequence

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} \frac{P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)}{P(w_1 \dots w_n)} \text{ Bayes Rule}$$

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n) \text{ Dropping the denominator}$$

121

Decoding with Hidden Markov Models (ii)

$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} \frac{\text{"Likelihood"}^n}{\text{"Prior"}^n}$

- HMM taggers make two further simplifying assumptions.
 - The probability of a word appearing depends only on its own tag and is independent of neighboring words and tags: $P(w_1 \dots w_n | t_1 \dots t_n) \approx \prod_{i=1}^n P(w_i | t_i)$
 - The second assumption, the bigram assumption, is that the probability of a tag is dependent only on the previous tag, rather than the entire tag sequence; $P(t_1 \dots t_n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$

122

Decoding with Hidden Markov Models (iii)

- Plugging the simplifying assumptions results in the following equation for the most probable tag sequence from a bigram tagger:

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n) \approx \operatorname{argmax}_{t_1 \dots t_n} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission transition}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$$

- The two parts correspond neatly to the **B** emission probability and **A** transition probability that we defined previously!

123

SIT330-770: Natural Language Processing

Week 5.14 – The Viterbi Algorithm

Dr. Mohamed Reda Bouadjenek

School of Information Technology,
Faculty of Sci Eng & Built Env

124

Computing the most probable sequence of tags

A brute force approach to identify the most probable sequence of tags faces exponential complexity

- This method is impractical for large datasets or real-time applications.

Solution: The Viterbi algorithm 1967

- Leverages dynamic programming, streamlining the process by breaking the problem into manageable sub-problems
- This approach significantly reduces computational demands and enhances processing speed, making it viable for complex tasks in real-world scenarios

Andrew Viterbi

125

The Viterbi Algorithm (i)

- The decoding algorithm for HMMs is the **Viterbi algorithm**
 - As an instance of dynamic programming, Viterbi resembles the dynamic programming minimum edit distance algorithm
 - The Viterbi algorithm first sets up a probability matrix or lattice:
 - Columns as **observables** (words of a sentence in the same sequence as in sentence)
 - Rows as hidden states (all possible POS Tags are known)

tag the sentence
Janet will back the bill

126

The Viterbi Algorithm (ii)

- Each cell of the matrix is represented by $V_{t,j}$ (Viterbi value for t: column, j: row) having the probability that the HMM is in **state j** (present POS Tag) after seeing the **first observations** (past words for which matrix (cell) values has been calculated) and passing through the most **probable state sequence (previous POS Tag)** q_1, \dots, q_{t-1}
- Computed by recursively taking the most probable path that could lead us to this cell

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

$v_{t-1}(i)$ the previous Viterbi path probability from the previous time step
 a_{ij} the transition probability from previous state q_i to current state q_j
 $b_j(o_t)$ the state observation likelihood of the observation symbol o_t given the current state j

127

The Viterbi Algorithm (iii)

- Each cell of the matrix is represented by $V_{t,j}$ (Viterbi value for t: column, j: row) having the probability that the HMM is in **state j** (present POS Tag) after seeing the **first t observations** (past words for which matrix (cell) values has been calculated) and passing through the most **probable state sequence (previous POS Tag)** q_1, \dots, q_{t-1}

A sketch of the matrix for Janet will back the bill, showing the possible tags (q) for each word and highlighting the path corresponding to the correct tag sequence through the hidden states

States (parts of speech) which have a zero probability of generating a particular word according to the B matrix (such as the probability that a determiner DT will be realized as Janet) are greyed out

128

Working Example (i)

- Janet will back the bill \rightarrow Janet/NNP will/MD back/VB the/DT bill/NN

	NNP	MD	VB	JJ	NN	RB	DT
S<>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
NNP	0.3777	0.0110	0.0099	0.0084	0.0084	0.0090	0.0225
MD	0.0006	0.0006	0.7908	0.0005	0.0005	0.0004	0.0001
VB	0.0322	0.0006	0.0006	0.0005	0.0015	0.0514	0.0231
JJ	0.0366	0.0004	0.0001	0.0733	0.0459	0.0035	0.0036
NN	0.0096	0.0178	0.0014	0.0086	0.1216	0.0177	0.0068
RB	0.0084	0.0006	0.0006	0.0006	0.0006	0.0006	0.0039
DT	0.1147	0.0021	0.0002	0.2157	0.0454	0.0102	0.0017

The A transition probabilities $P(q_t|q_{t-1})$ computed from the WSJ corpus without smoothing

Observation likelihoods B computed from the WSJ corpus without smoothing, simplified slightly

129

Working Example (ii)

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t)$$

$v_{t-1}(i)$ the previous Viterbi path probability from the previous time step
 a_{ij} the transition probability from previous state q_i to current state q_j
 $b_j(o_t)$ the state observation likelihood of the observation symbol o_t given the current state j

130

Evaluation Metrics for Named Entity Recognition (NER)

- NER Evaluation Basics:**
 - Unlike POS tagging, evaluated on accuracy, NER uses recall, precision, and F1 score.
 - Recall measures correctly identified entities against all actual entities.
 - Precision counts correct labels against all labeling attempts.
 - The F1 score provides a balance between precision and recall, serving as a single metric for accuracy.
- Challenges in NER:**
 - NER systems treat entities as single units for evaluation, leading to challenges not seen in POS tagging.
 - The system's ability to correctly identify entire entities, such as 'Jane Villanueva', impacts evaluation outcomes.
 - Mismatches in entity recognition across training and test data can skew results.

132