

Natural Language Processing

Week 6

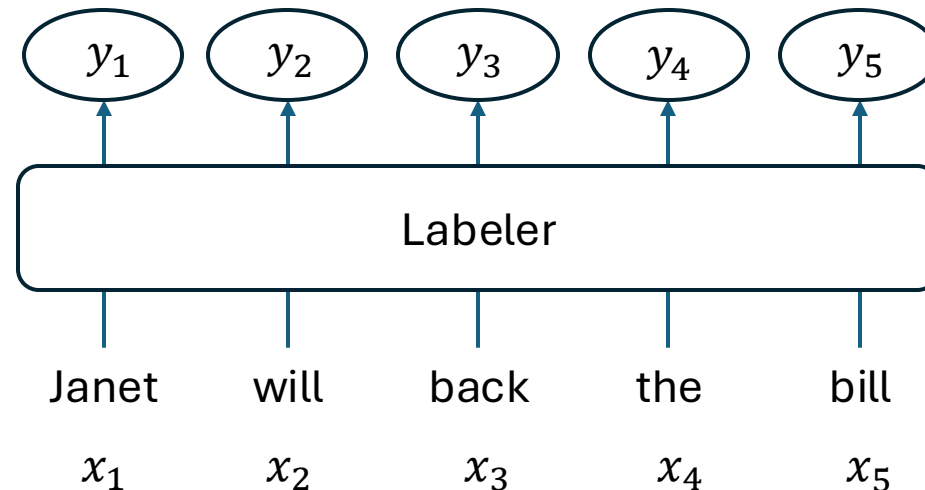
Agenda

- Project Introduction
- Token Classification
 - POS
 - NER
- Extractive Summarization

Token Classification

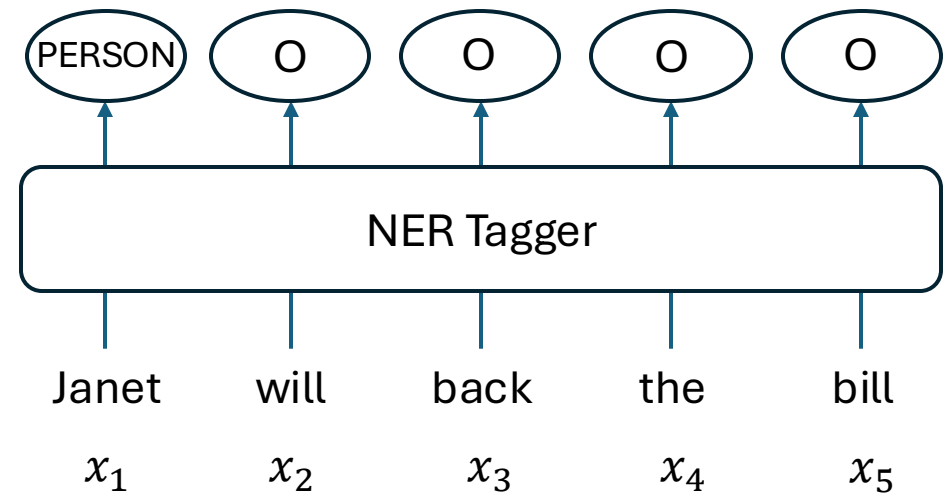
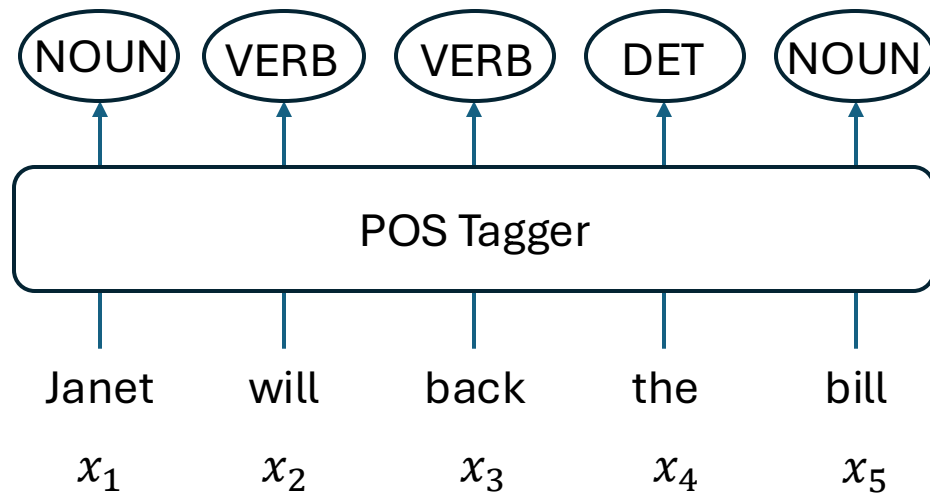
Token Classification

- Token classification is the task of assigning a categorical label for every element in sequence of data
 - Text data, biological sequences, time-series
- For NLP applications, the task is to **assign each word x_i in an input sequence a label y_i** (also called **sequence labeling**)



Token Classification

- **Parts-of-Speech (POS) tagging** and **Named Entity Recognition (NER)** are two of the most common forms of sequence modeling in NLP



Parts of Speech

Parts of Speech

- Parts of speech are formally defined based on their grammatical relationship with neighboring words or the morphological properties of their affixes

	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>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by, under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, 2026, 11:00, hundred</i>
	PART	Particle: a function word that must be associated with another word	<i>'s, not, (infinitive) to</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
Other	SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>whether, because</i>
	PUNCT	Punctuation	<i>; , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

Parts of Speech

- **Open class:** category of parts of speech that is open to being expanded / updated / changed as language evolves
 - Nouns: bitcoin, iPhone, etc
 - [We've updated the Merriam-Webster.com Dictionary with 690 New Words | Merriam-Webster](#)
- **Closed class:** category of parts of speech that is generally fixed / doesn't evolve
 - Prepositions: of, about, around

Parts of Speech (Penn Treebank)

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	infinitive to	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>’s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

Evolution of POS Tagging

- **Rule-Based Tagging (60s-80s):** manual tagging and rule-based systems utilizing dictionaries and manually engineering grammatical rules
- **Probabilistic Methods (80s-90s):** Hidden Markov Models (HMMs) tagged sequences based on word occurrence probabilities, leading to better ambiguity resolution than rule-based systems
- **Machine Learning Methods (2000s):** Conditional Random Fields (CRFs) and other ML algorithms utilized feature context and interactions
- **Deep Learning Methods (2010s):** Bidirectional RNNs and LSTMs enhanced handling of long-range dependencies, minimizing manual feature engineering.
- **Transformers (current):** Efficiency, multi-language improvements, etc.

Task Performance

- Accuracy of part-of-speech tags is very high
- 97% accuracy across 15 languages from Universal Dependency treebank
- Accuracies for other treebanks are 97% no matter the algorithm (Hidden Markov Model, Conditional Random Field, BERT, etc)

Parts of Speech: How hard is the task?

- Most word types (unique words) are unambiguous (85-86%)
- Ambiguous tokens (instances of words), however, appear very often
- Most-frequent-tag baseline gives 92% accuracy (STOA is 97%)

Types:		WSJ	Brown
Unambiguous	(1 tag)	44,432 (86%)	45,799 (85%)
Ambiguous	(2+ tags)	7,025 (14%)	8,050 (15%)
Tokens:			
Unambiguous	(1 tag)	577,421 (45%)	384,349 (33%)
Ambiguous	(2+ tags)	711,780 (55%)	786,646 (67%)

earnings growth took a **back/JJ** seat
a small building in the **back/NN**
a clear majority of senators **back/VBP** the
bill

Dave began to **back/VB** toward the door
enable the country to buy **back/RP** debt
I was twenty-one **back/RB** then

Parts of Speech: Applications

- POS tagging could be used for:
 - NER Models (creating features, used for post-processing)
 - Models for sentence segmentation (features)
 - Sentiment analysis (features)
 - Improving speech recognition
- Currently, neural networks and modern transformer architecture **has reduced the need to explicitly engineer features for NLP applications**
- However, there may still be places where POS is needed (environments that can't support transformers models)

Named Entity Recognition

Named Entity Recognition (NER)

- Named entities are words / phrases in text that refer to proper nouns
- NER is an NLP task in which the goal is to find **spans of text** that constitute proper names and **assign the correct entity type** to the span
- Most common are **people, places, organizations**, or **geo-political** entities
- Does not have to be an entity per se (could be time, or any other short span of text you deem as an entity (such as cost of a contractual agreement))

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].

NER Labels

- For NER, every word in a sequence is labeled
- **BIO/IOB2** – label token that begins a span with **B**, tokens that occur inside are labeled with **I**, tokens outside the spans are with **O**
- **IO** – loses indicator for the beginning of a span
- **BIOES/IOBES** – adds indicator for end of a span, and a single-token span

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

Evaluation of NER

- NER is a highly **imbalanced, multi-class classification task**
- Precision, Recall, F-score as standard evaluation metrics
- NER has two dimensions of correctness: **type** and **span**
- A popular library call **seqeval** penalizes on both dimensions
- Other libraries exist that may not be as strict: [nervaluate · PyPI](#)

Evaluation of NER

tokens	true	underspan	overspan	type
Steve	B-PER	B-PER	B-PER	B-PER
Jobs	I-PER	O	I-PER	I-PER
and	O	O	O	O
Tim	B-PER	B-PER	B-PER	B-PER
Cook	I-PER	O	I-PER	I-PER
both	O	O	O	O
led	O	O	O	O
Apple	B-ORG	B-ORG	B-ORG	B-LOC
,	O	O	I-ORG	O
which	O	O	I-ORG	O
is	O	O	O	O
based	O	O	O	O
in	O	O	O	O
Cupertino	B-LOC	B-LOC	B-LOC	B-LOC
.	O	O	O	O

UNDER

OVER

TYPE

Accuracy: 0.87

	precision	recall	f1-score	support
LOC	1.00	1.00	1.00	1
ORG	1.00	1.00	1.00	1
PER	0.00	0.00	0.00	2
micro avg	0.50	0.50	0.50	4
macro avg	0.67	0.67	0.67	4
weighted avg	0.50	0.50	0.50	4

Accuracy: 0.87

	precision	recall	f1-score	support
LOC	1.00	1.00	1.00	1
ORG	0.00	0.00	0.00	1
PER	1.00	1.00	1.00	2
micro avg	0.75	0.75	0.75	4
macro avg	0.67	0.67	0.67	4
weighted avg	0.75	0.75	0.75	4

Accuracy: 0.93

	precision	recall	f1-score	support
LOC	0.50	1.00	0.67	1
ORG	0.00	0.00	0.00	1
PER	1.00	1.00	1.00	2
micro avg	0.75	0.75	0.75	4
macro avg	0.50	0.67	0.56	4
weighted avg	0.62	0.75	0.67	4

Span-Flexible Evaluation of NER

tokens	true	underspan	overspan	type
Steve	B-PER	B-PER	B-PER	B-PER
Jobs	I-PER	O	I-PER	I-PER
and	O	O	O	O
Tim	B-PER	B-PER	B-PER	B-PER
Cook	I-PER	O	I-PER	I-PER
both	O	O	O	O
led	O	O	O	O
Apple	B-ORG	B-ORG	B-ORG	B-LOC
,	O	O	I-ORG	O
which	O	O	I-ORG	O
is	O	O	O	O
based	O	O	O	O
in	O	O	O	O
Cupertino	B-LOC	B-LOC	B-LOC	B-LOC
.	O	O	O	O

UNDER

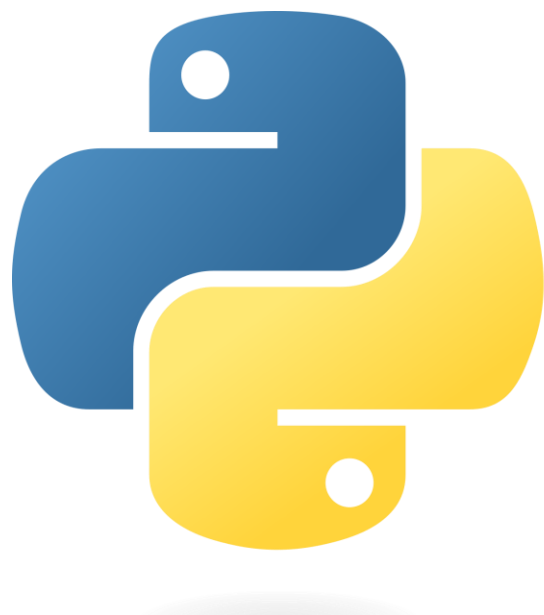
	precision	recall	f1
LOC	1.0	1.0	1.0
ORG	1.0	1.0	1.0
PER	1.0	1.0	1.0

OVER

	precision	recall	f1
LOC	1.0	1.0	1.0
ORG	1.0	1.0	1.0
PER	1.0	1.0	1.0

TYPE

	precision	recall	f1
LOC	1.0	1.0	1.0
ORG	0.0	0.0	0.0
PER	1.0	1.0	1.0



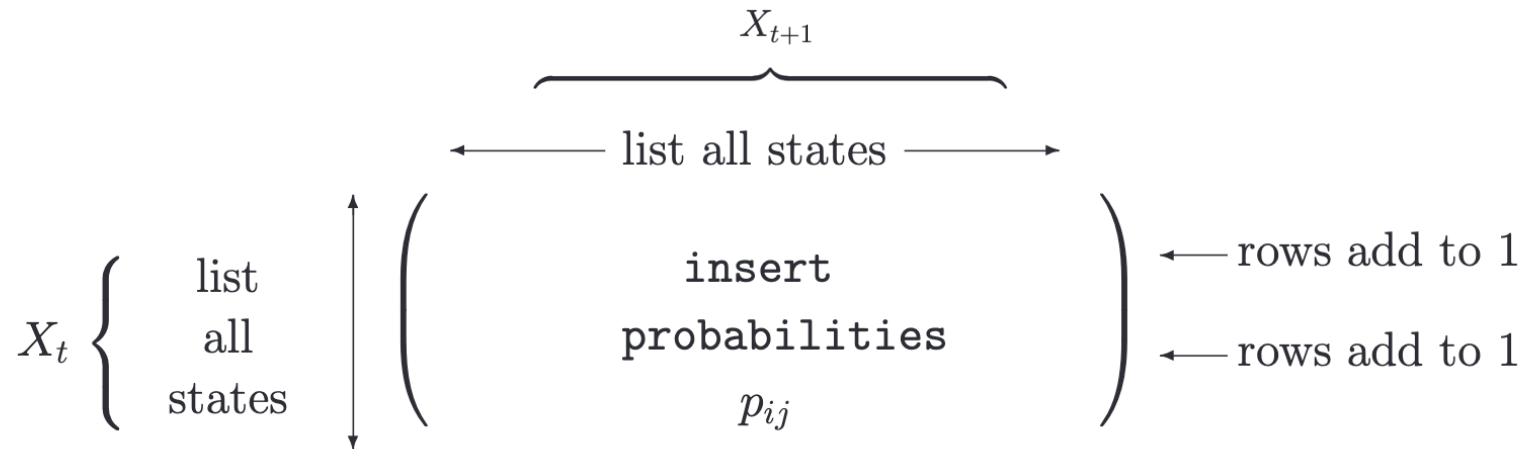
Extractive Summarization

Summarization

- There are two types of summarization:
 - **Extractive:** Method wherein key sentences are **extracted** without modification from the original document.
 - **Abstractive:** Method wherein new sentences are **generated** to represent a summary of the original document.
- Extractive Pro:
 - Faithful to the original text
- Extractive Con:
 - Can lead summaries that are not fluent

PageRank

- The text rank algorithm is based of the **PageRank** algorithm
- Given a matrix of transition probabilities for web pages (web pages as nodes, links between them as edges)



- Use PageRank algorithm to compute stationary distribution (long-term probabilities)

$$r_{new} = (\alpha P + (1 - \alpha)E) \times r_{old}$$

- Higher values in PageRank indicate a webpage is “more important”

TextRank

- Instead of a transition matrix of webpages, **TextRank relies on similarities between sentences** (sentences as nodes, similarities as edges)
- Sentences are represented as vectors
 - Embeddings
 - BOW
- The pairwise similarity of those vectors is calculated
 - Cosine Similarity
 - Jaccard
- The core convergence algorithm is still applied
- **A sentence is considered important if it is similar to many other sentences**

