
Computational Linguistics for Pakistan Languages

Tafseer Ahmed

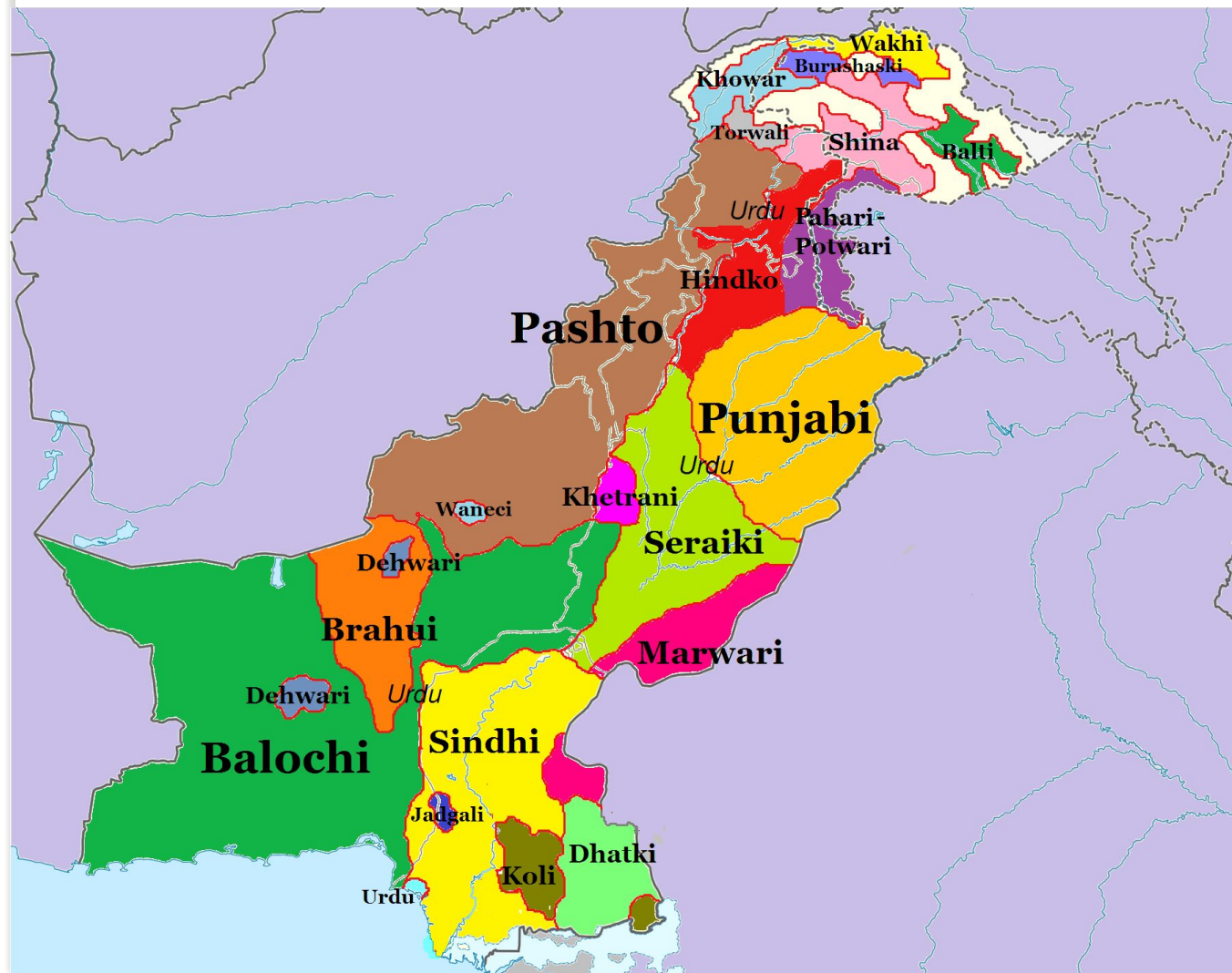
Languages in Pakistan

- **Languages: 73**, Indigenous: 65
- Institutional: 7, Developing: 17, Vigorous: 39, Trouble: 8, Dying: 2
- Families: Indo-Aryan, Iranian, Dravidian, Tibetan, Isolate
- <https://www.ethnologue.com/country/PK>



Overview





Linguistic Analysis

- **Traditional Grammars**

- paragraphs explaining the grammar

- **Formal Grammars**

- Formal grammar rules
- Panini (4th century BC) , Chomsky (1928-)

Computer Science

- **Computational Linguistics (CL)**
 - modeling of natural (human) languages
- **Natural Language Processing (NLP)**
 - processing of natural languages

Different Techniques of NLP

- Linguistic Analysis
- Non-Linguistic Methods
 - Bag of Words
 - Sequence to Sequence

Linguistic Analysis

- Deep Method *(not related to deep learning)*
 - Inspired by Linguistics
-
- Token → Part of Speech → Morphology → Syntax → Semantics → Discourse → ...

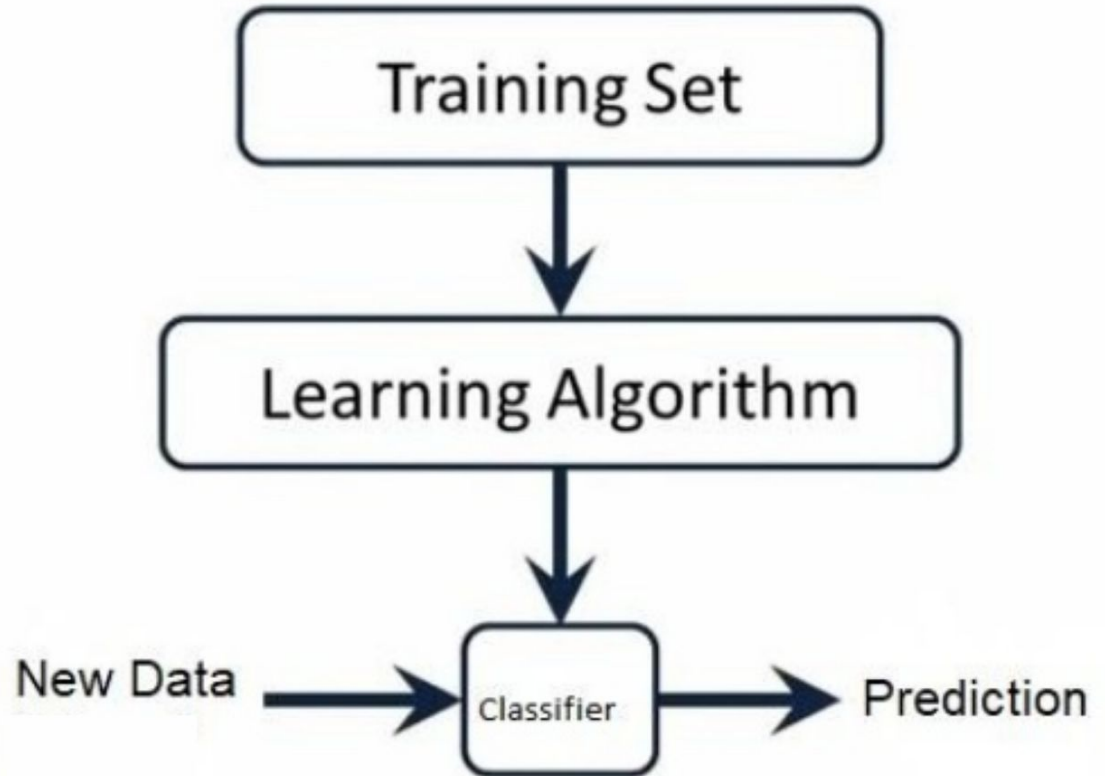
Linguistic Features

- Linguistic Features are used as features in many machine learning systems.
- Many current and old systems uses linguistic analysis, however modern trend is End-to-End systems.

Non-Linguistic Methods

- **Bag of Words** - Order does not matter
 - Text Classification, Text Clustering, Topic Modeling, ...
- **Sequence Modeling** - Order does matter
 - Language Modeling, Machine Translation, ...

Supervised Machine Learning



A sample feature set in SVMTool

word features	$w_{-3}, w_{-2}, w_{-1}, w_0, w_{+1}, w_{+2}, w_{+3}$
PoS features	$p_{-3}, p_{-2}, p_{-1}, p_0, p_{+1}, p_{+2}, p_{+3}$
ambiguity classes	a_0, a_1, a_2, a_3
may_be's	m_0, m_1, m_2, m_3
word bigrams	$(w_{-2}, w_{-1}), (w_{-1}, w_{+1}), (w_{-1}, w_0)$ $(w_0, w_{+1}), (w_{+1}, w_{+2})$
PoS bigrams	$(p_{-2}, p_{-1}), (p_{-1}, a_{+1}), (a_{+1}, a_{+2})$
word trigrams	$(w_{-2}, w_{-1}, w_0), (w_{-2}, w_{-1}, w_{+1}),$ $(w_{-1}, w_0, w_{+1}), (w_{-1}, w_{+1}, w_{+2}),$ (w_0, w_{+1}, w_{+2})
PoS trigrams	$(p_{-2}, p_{-1}, a_{+0}), (p_{-2}, p_{-1}, a_{+1}),$ $(p_{-1}, a_0, a_{+1}), (p_{-1}, a_{+1}, a_{+2})$
sentence_info	punctuation ('.', '??', '!')
prefixes	$s_1, s_1s_2, s_1s_2s_3, s_1s_2s_3s_4$
suffixes	$s_n, s_{n-1}s_n, s_{n-2}s_{n-1}s_n, s_{n-3}s_{n-2}s_{n-1}s_n$
binary word features	initial Upper Case, all Upper Case, no initial Capital Letter(s), all Lower Case, contains a (period / number / hyphen ...)
word length	integer

Linguistic Analysis - Goals (of this series)

- Learning and creating **Linguistic Representations** that
 - are widely used in computer applications
 - deals with peculiar features of Pakistani languages
- Training the **annotators** for creating linguistically correct training dataset

Linguistic Analysis - Goals (of this series)

- Creating linguistic analyzers by using **small training datasets**
 - Unsupervised Learning
 - Word Embeddings
 - Transfer Learning
 -
- Using the **libraries and tools** of computational linguistics

Not the Goals/Focus of this series

- Text Mining
- Corpus Linguistics
- NLP algorithms - in depth
- Machine (including Deep) Learning - in depth
 - However, the topics related to linguistic analysis will be discussed

Linguistic Analysis

Part of Speech (PoS)

پڑھی	کتاب	اچھی	ایک	روزانہ	نے	لڑکی	ذہین
VERB	NOUN	ADJ	NUM	ADV	ADP	NOUN	ADJ

PoS - Universal Tagset

Open class words	Closed class words	Other
<u>ADJ</u>	<u>ADP</u>	<u>PUNCT</u>
<u>ADV</u>	<u>AUX</u>	<u>SYM</u>
<u>INTJ</u>	<u>CCONJ</u>	<u>X</u>
<u>NOUN</u>	<u>DET</u>	
<u>PROPN</u>	<u>NUM</u>	
<u>VERB</u>	<u>PART</u>	
	<u>PRON</u>	
	<u>SCONJ</u>	

Building PoS Tagger using small dataset

- The available PoS Tagger(s) for Urdu are not very good.
- Annotated Data is required for other Pakistani languages

Morphological Features

تھیں	پڑھی	کتابیں	اچھی	نے	لڑکیوں	ذہین	
ہے	پڑھ	کتاب	اچھا	نے	لڑکی	ذہین	Lemma
AUX	VERB	NOUN	ADJ	ADP	NOUN	ADJ	POS
	Form=Perf Gend=Fem Pers=3	Gend=Fem Num=Pl Form=Nom	Gend=Fem Num=Pl Form=		Gend=Fem Num=Pl Form=Obl	Gend=Fem Num=Pl Form=Obl	MORPH. Features

Morphology - Pakistani Languages

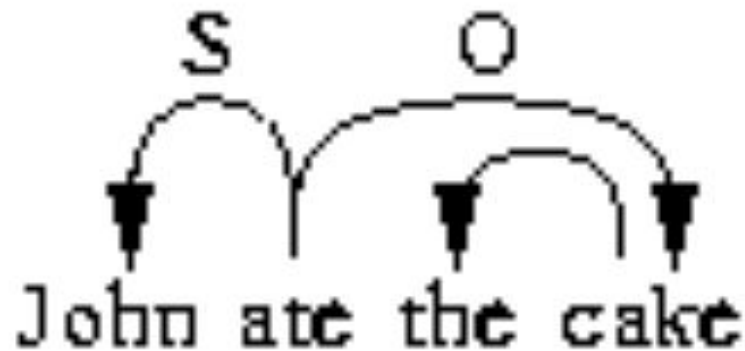
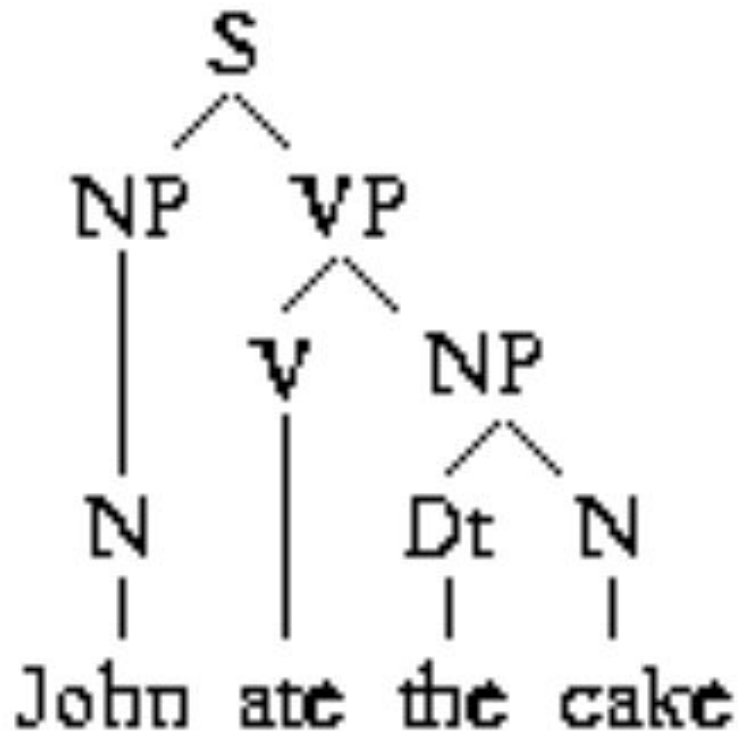
- **Rich Morphology**
- Nominative and Oblique Forms
- Inflectional Case
- Causatives
- Reduplication
- Pronominal Suffixes

Derivational Morphology

Morphological Analysis

- Finite State Transducer
- Morph. Feature Prediction
- Unsupervised Learning and Prediction

Syntax: Phrase Structure vs Dependencies



Universal Dependencies

Core dependents of clausal predicates

Nominal dep	Predicate dep	
<u>nsubj</u>	<u>csbj</u>	
<u>nsubjpass</u>	<u>csbjpass</u>	
<u>dobj</u>	<u>ccomp</u>	<u>xcomp</u>
<u>iobj</u>		

Noun dependents

Nominal dep	Predicate dep	Modifier word
<u>nummod</u>	<u>acl</u>	<u>amod</u>
<u>appos</u>		<u>det</u>
<u>nmod</u>		<u>neg</u>

Non-core dependents of clausal predicates

Nominal dep	Predicate dep	Modifier word
<u>nmod</u>	<u>advcl</u>	<u>advmod</u>
		<u>neg</u>

Compounding and unanalyzed

<u>compound</u>	<u>mwe</u>	<u>goeswith</u>
<u>name</u>	<u>foreign</u>	

Special clausal dependents

Nominal dep	Auxiliary	Other
<u>vocative</u>	<u>aux</u>	<u>mark</u>
<u>discourse</u>	<u>auxpass</u>	<u>punct</u>
<u>expl</u>	<u>cop</u>	

Coordination

<u>conj</u>	<u>cc</u>	<u>punct</u>
-------------	-----------	--------------

تھیں	پڑھیں	کتابیں	اچھی	نے	لڑکیوں	ذہین	
ہے	پڑھ	کتاب	اچھا	نے	لڑکی	ذہین	Lemma
Aux	Verb	NN	Adj	AdP	Noun	Adj	POS
	Form=Perf Gend=Fem Pers=3 Num=Pl	Gend=Fem Num=Pl Form=Nom	Gend=Fem Num=Pl Form=Nom		Gend=Fem Num=Pl Form=Obl	Gend=Fem Num=Pl Form=Obl	Features

Universal Dependencies - To Do

- Revisiting annotation scheme for Urdu (and other South Asian languages)
- Creating UD_banks for other Pakistani languages

Annotation

- **GATE** (General Architecture for Text Engineering), Sheffield
- **Brat**, Manchester, Tokyo,.....
- **Webanno**, Darmstadt
-

dependency-conll - Notepad									
File	Edit	Format	View	Help					
1	ذبین	ذبین	Adj	Adj	-	2	amod		
2	لڑکیاں	لڑکی	Noun	NN	-	6	subj		
3	نے	نے	Adp	PP	-	2	case		
4	اچھی	اچھا	Adj	Adj	-	5	amod		
5	کتابیں	کتاب	Noun	NN	-	6	obj		
6	پڑھیں	پڑھ	Verb	VB	-	0	ROOT		
7	تھیں	ہے	Aux	Aux	-	6	aux		

Semantic Role

- The **glass**[patient] broke.
- **I**[agent] break the **glass**[patient].
- The **glass**[patient] was broken by **me**[agent].

Urdu PropBank

Predicate: *DagmagA*

DagmagA: created by Tafseer - Hindi dagmagA

Roleset id: **DagmagA.01** , *To stagger/tremble/shake*

Roles:

Arg1: *The entity that staggers/trembles/shakes*

Example:

<http://www.express.pk/story/44282/>

بیلی کایٹر بُری طرح ڈگمگایا

Arg1: بیلی کایٹر

Argm-mnr: بُری طرح

Rel: ڈگمگایا

Lexical Resources

- Urdu Wordnet
- <http://www.cle.org.pk/clestore/urduwordnet.htm>

اُٹ پٹ	ان پٹ
Noun 100034 : کسی بھی ملک یا خطے کی زبان آپ کی لسان عربی ہو یا اردو رہیں گے آپ مسلمان ہی بھاشا، بولی، زبان، لسان	78172
Noun 102523 : منہ کے اندر کا وہ عضو جس میں قوت ذائقہ ہوتی ہے اور جو نطق کا آلہ ہے گرم چائے سے اُس کی زبان پر چھالے پڑ گئے جیبھ، زبان، لسان	دکھائی دینے والا متن داخل کریں
Noun 104948 : قول، کہی ہوئی بات، وعدہ کچھ بھی ہو جائے میں زبان دے کر کبھی نہیں مکتا زبان	خالی کریں تخصیص کریں

Distributional Semantics

- “You shall know a word by the company it keeps.” John Rupert Firth
- Word Embeddings
- Topic Models
-



```

: model.most_similar('سورج', topn=10)

: [ ('0.8165232539176941', 'چاند'),
  ('0.814332127571106', 'سُورج'),
  ('0.7290103435516357', 'آفتاب'),
  ('0.7133749723434448', 'ماہتاب'),
  ('0.707730770111084', 'تارے'),
  ('0.702995777130127', 'ستاروں'),
  ('0.69676274061203', 'ستارے'),
  ('0.6915234327316284', 'کرنوں'),
  ('0.6877450942993164', 'مریخ'),
  ('0.6747698783874512', 'غروب'),
  ('0.6725152730941772', 'بادلوں'),
  ('0.668827474117279', 'تاروں'),
  ('0.6525486707687378', 'گہن'),
  ('0.6493672132492065', 'شفق'),
  ('0.6463797092437744', 'شعاعوں') ]

```


Discourse Analysis

- He wrote that last year. ...

Who is HE?

What is THAT?

When is the LAST YEAR?

Filling the gaps – Pro Drop

- وہ کمرے میں آیا، کرسی ڈھونڈی اور بیٹھ گیا
- وہ کمرے میں آیا، (اس نے) کرسی ڈھونڈی اور (وہ) (اس پر) بیٹھ گیا

Thank You

شکریہ

Contact: tafseer@gmail.com