# Learning-Based Named Entity Recognition for Morphologically-Rich, Resource-Scarce Languages

Kazi Saidul Hasan Md. Altaf ur Rahman Vincent Ng

Human Language Technology Research Institute
University of Texas at Dallas

European Chapter of the Association for Computational Linguistics, 2009

# Goal

Named Entity Recognition (NER) for morphologically-rich, resource-scarce languages

- Use Bengali as our representative language
- Focus on identifying PERSON, ORGANIZATION, and LOCATION

# Challenges

- Scarcity of NE-labeled data
  - Unsupervised techniques for English NER (e.g., Collins and Singer (1999)) are unlikely to work well
    - Problem: Lack of capitalization in Bengali
- Lack of publicly available gazetteers
- Inaccurate POS tagger
  - Unsupervised POS induction techniques (e.g., distributional clustering) are unlikely to work well
     Problems: Distributional representation may not be reliable for Bengali because of
    - free word order
    - morphological inflections

# Weaknesses of Existing Bengali NE Recognizers

- Use their own manually-constructed gazetteers
  - The results are not reproducible
- Use pseudo-affixes (created by extracting the first n and the last n characters of a word)
  - The process is ad-hoc and may not cover many useful affixes
- Typically adopt a pipelined NER architecture
  - Errors from the POS tagger are propagated to the NE recognizer

# Our Approach

## Still supervised, but

- we investigate two new features when used in a pipelined architecture for POS tagging and NER:
  - affixes induced from an unannotated corpus
  - semantic class information extracted from Wikipedia
- we propose a joint model for learning POS tagging and NER simultaneously

# Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- Pipelined NER with New Features
  - Evaluation
- 4 Joint Model for POS Tagging and NER
  - Evaluation

# Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- Pipelined NER with New Features
  - Evaluation
- Joint Model for POS Tagging and NER
  - Evaluation

# **Affix Induction**

## Goal

For morphologically-rich languages, a lot of grammatical information is expressed via affixes. Unlike previous approaches, we learn affixes rather than using pseudo-affixes.

# Approach: Keshava and Pitler's (2006) method

Assume that

- (1) *V* is a vocabulary extracted from a large, unannotated corpus
- (2)  $\alpha$  and  $\beta$  are two character sequences
- (3)  $\alpha\beta$  is the concatenation of  $\alpha$  and  $\beta$ 
  - if  $\alpha\beta$ ,  $\alpha \in V$ , we extract  $\beta$  as a candidate suffix
  - if  $\alpha\beta$ ,  $\beta \in V$ , we extract  $\alpha$  as a candidate prefix

# Affix Induction (contd.)

Many of the induced affixes could be spurious.

Example: If both "can" and "candidate" are in V, then "didate" is extracted as an induced suffix

## Solution

- Assign a score to each induced affix
- Select only those that score over a certain threshold

We use induced affix n-grams as features. These are like word n-grams except that these are made of affixes.

# Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- 3 Pipelined NER with New Features
  - Evaluation
- Joint Model for POS Tagging and NER
  - Evaluation

# Semantic Class Induction from Wikipedia

## Goal

- Generate a list of phrases and tokens that are potentially named entities from all the articles in the Bengali Wikipedia
- Peuristically annotate each of them with one of four classes: PER (person), ORG (organization), LOC (location), or OTHERS

# Generating an Annotated List of Phrases

# Steps

- Generate and annotate the title of each article
  - (i) Use category information to annotate the title



আনক অর্জন পিখ্যাগার্র্যাসর ধারণারই বিশ্রনি।

वकि सक विश्वकाष পরিত্রমন

প্রধান পাজা

সম্পদায

 সাহায্য ল দাল ককল

অনুসন্ধান

6(ना

সম্পর্কিত পরিবর্তন

ছাগার যোগ্য সংস্করণ

বিশেষ পৃষ্ঠাসমূহ

 স্থানী সংযোগ এই নিবন্ধটি উদ্ধৃত

অন্যান্য ভাষাসময়

 Asturianu Azərbaycan

v. Žemaitėška

Беларуская

Български

रेमात ठीत/विक्षिया

কক্ৰন

العربية =

হাজিয়ার সংযোগকারী পর্তাসমৃহ

वाश्ना উইकिशिक्षिया

 সমসাম্বিক ঘটনা সাম্প্রতিক পরিবর্তনসমহ

অজানা খেকোনো পর্তা

অনসন্ধান

উইকিপিডিয়া, মক্ত বিশ্বকোষ খেকে

শিশাগোরাস বা সামোসের শিশাগোরাস (প্রাচীন ত্রিক ভাষায় Ιωθαγόρας *পুখাগোরাস*) ( ত্রিস্টপূর্ব ৫৮০ - ত্রিস্টপূর্ব ৫০০) একজন গ্রীক দাশনিক ও গণিজবিদ।

পিখাগোরাস সবচেয়ে বেশি বিখ্যাভ পিখাগোরাসের উপপাদোর জন্য। গ্রীপ্টের জন্মের প্রায় ৬০০ বছর আগে পিখাগোরাস দর্শণ শাস্ত্রে প্রভুভ অবদান রাখেন। ভাকে বলা হয় - সংখ্যার জনক। পিখাগোরাস প্রথম নিজেকে দার্শনিক হিসাবে দাবী করেন। প্লেটা, এরিপটল ও কোপারনিকাসের

ভার কোনো দেখা পরে আর পাওয়া যায়নি ফলে ভার অবদান সম্পর্কে নিশ্চিত হওয়া যায়নি।

৫২৯ খ্রিস্টপর্বাদে ভিনি ইভাদীর ক্রোটোনা শহরে বসবাস শুরু করেন।

পিখাগোরাসকে বলা যায় সংখ্যা দার্শনিক। তিনি বিশ্বাস করতেন পৃথিবীর সকল কিছর মল হলো সংখ্যা। সংখ্যার সঙ্গে নৈতিকতা, রঙ, শুভ-অশুভকে তিনি মিরিয়ে ফেলেন। তার বিশ্বাস হলো মানুষের আন্না অবিনশ্বর। ফলে হিন্দু, বৌদ্ধ ও অন্যান্য ক্ষেকটি ধর্মের মতো তিনিও জন্মান্তরবাদে বিশ্বাসী। পিখাগোরাস পথিবীর গোলকত্বে বিশ্বাস করতেন কারণ ভার ধারণা ছিল বত হলো পরম বক্ররেখা। সর্য, চন্দ্র বা গ্রহদের



#### পিখাগোরাদের প্রতিকৃতি

[সম্পাদনা]

পিখাগোরাসের উপপাদ্য সমকোণী ত্রিভ্তের অভিভ্তের ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফল অপর দুই বাহুর ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফলের সমষ্টির সমান।

নিজেদের একটা গতি আছে।

এই निवन्निक्षि व्यमस्पूर्ण। व्यापनि हारेल अविक ममृद्ध कराल भारतन।

7 - ST - 7

প্লাক-সক্রেটীয় দার্শনিকবন্দ

মিলেভসীয় ধারা: (খলিস আনাক্সিমান্ড্রোস মিলেভুসের আনাক্সিমেনিস **পিখাপোরাসবাদী: পিখাপোরাস** • ফিলোলাউস • আন্ধরামোন • আর্থভাস • তিমামোস

এফোস্সীয় ধারা: ইরাইইডাস - এল্যার দার্শনিক ধারা: জোনাফানিজ - গামনিদস - এল্যার জিনা - মেলিসাস বহুত্বাদী ধারা: আন্মোগোরাস · এম্পোদারুস — পরমাণবাদী ধারা: দেউকিয়োস · দেমাক্রিভোস

সোহিবাদ: প্রোজাগারাস · গোর্গিবাস · প্রোদিকোস · ইপ্লিবাস

দিওগেনেস - ফেরেকদেস

বিষয়প্রেণীসমহ: অসম্পর্ণ । গ্রিক গণিতবিদ । গ্রিক দার্শনিক । খ্রিস্টপর্ব ৫৮০-এ জন্ম । খ্রিস্টপর্ব ৫০০-এ মত্য

মণিপুরী Brezhonea Bosanski



# **Pythagoras**

## পরিরমন

- প্রধান পাজা
- वाश्ना উইकिशिक्षिया সম্পদায

वकि सक विश्वकाष

- সমসাম্বিক ঘটনা সাম্প্রতিক পরিবর্তনসমহ
- অজানা খেকোনো পর্তা
- সাহায্য
- ল দাল ককল

## অনুসন্ধান

6(ना



## হাজিয়ার

- সংযোগকারী পর্তাসমৃহ সম্পর্কিত পরিবর্তন
- বিশেষ পৃষ্ঠাসমূহ ছাগার যোগ্য সংস্করণ
- স্থানী সংযোগ
- এই নিবন্ধটি উদ্ধৃত কক্ৰন

#### অন্যান্য ভাষাসময়

- العربية = Asturianu
- Azərbaycan v. Žemaitėška
- Беларуская
- Български
- Brezhonea Bosanski

 रेमात ठीत/विक्षिया মণিপুরী

শিশাগোরাস বা সামোসের শিশাগোরাস (প্রাচীন ত্রিক ভাষায় Ιωθαγόρας *পুখাগোরাস*) ( ত্রিস্টপূর্ব ৫৮০ - ত্রিস্টপূর্ব ৫০০) একজন গ্রীক দাশনিক ও গণিজবিদ।

পিখাগোরাস সবচেয়ে বেশি বিখ্যাভ পিখাগোরাসের উপপাদোর জন্য। গ্রীপ্টের জন্মের প্রায় ৬০০ বছর আগে পিখাগোরাস দর্শণ শাস্ত্রে প্রভুভ অবদান রাখেন। ভাকে বলা হয় - সংখ্যার জনক। পিখাগোরাস প্রথম নিজেকে দার্শনিক হিসাবে দাবী করেন। প্লেটা, এরিপটল ও কোপারনিকাসের আনক অর্জন পিখ্যাগার্র্যাসর ধারণারই বিশ্রনি।

ভার কোনো দেখা পরে আর পাওয়া যায়নি ফলে ভার অবদান সম্পর্কে নিশ্চিত হওয়া যায়নি।

৫২৯ খ্রিস্টপর্বাদে ভিনি ইভাদীর ক্রোটোনা শহরে বসবাস শুরু করেন।

পিখাগোরাসের উপপাদ্য

7 - ST - 7

পিখাগোরাসকে বলা যায় সংখ্যা দার্শনিক। তিনি বিশ্বাস করতেন পৃথিবীর সকল কিছর মল হলো সংখ্যা। সংখ্যার সঙ্গে নৈতিকতা, রঙ, শুভ-অশুভকে তিনি মিরিয়ে ফেলেন। তার বিশ্বাস হলো মানুষের আন্না অবিনশ্বর। ফলে হিন্দু, বৌদ্ধ ও অন্যান্য ক্ষেকটি ধর্মের মতো তিনিও জন্মান্তরবাদে বিশ্বাসী। পিখাগোরাস পথিবীর গোলকত্বে বিশ্বাস করতেন কারণ ভার ধারণা ছিল বত হলো পরম বক্ররেখা। সর্য, চন্দ্র বা গ্রহদের নিজেদের একটা গতি আছে।

সমকোণী ত্রিভ্তের অভিভ্তের ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফল অপর দুই বাহুর ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফলের সমষ্টির সমান।



[प्रक्लापना]

এই निवन्निष्टि अमन्यूर्ण। आधिन हारेल अहिक ममृद्ध कराछ थातन।

প্লাক-সক্রেটীয় দার্শনিকবন্দ

মিলেভসীয় ধারা: (খলিস আনাক্সিমান্ড্রোস মিলেভুসের আনাক্সিমেনিস **পিখাপোরাসবাদী: পিখাপোরাস** • ফিলোলাউস • আন্ধরামোন • আর্থভাস • তিমামোস

এফোস্সীয় ধারা: ইরাইইডাস - এল্যার দার্শনিক ধারা: জোনাফানিজ - গামনিদস - এল্যার জিনা - মেলিসাস

বহুত্বাদী ধারা: আন্মোগোরাস · এম্পোদারুস — পরমাণবাদী ধারা: দেউকিয়োস · দেমাক্রিভোস

সোকিবাদ: (প্রাভাগোরাস • গোর্গিয়াস • প্রোদিকোস • ইপ্রিযাস দিওগেনেস - ফেরেকদেস

বিষয়প্রেণীসমহ: অসম্পর্ণ । গ্রিক গণিতবিদ । গ্রিক দার্শনিক । খ্রিস্টপর্ব ৫৮০-এ জন্ম । খ্রিস্টপর্ব ৫০০-এ মত্য



## **Pythagoras**

## পরিরমন

- প্রধান পাজা
- সম্পদায
- সমসাম্বিক ঘটনা সাম্প্রতিক পরিবর্তনসমহ
- অজানা খেকোনো পর্তা

वकि सक विश्वकाष

- সাহায্য
- ল দাল ককল

## অনুসন্ধান

চলো



## সংযোগকারী পৃষ্ঠাসমৃহ

- সম্পর্কিত পরিবর্তন
- বিশেষ পৃষ্ঠাসমূহ ছাগার যোগ্য সংস্করণ
- স্থানী সংযোগ
- এই নিবন্ধটি উদ্ধৃত ককল
- অন্যান্য ভাষাসময়

- العربية = Asturianu
- Azərbaycan
- v. Žemaitėška
- Беларуская Български
- रेमात ठीत/विक्षिया
- মণিপরী Brezhonea Bosanski

শিশাগোরাস বা সামোসের শিশাগোরাস (প্রাচীন ত্রিক ভাষায় Ιωθαγόρας *পুখাগোরাস*) ( ত্রিস্টপূর্ব ৫৮০ - ত্রিস্টপূর্ব ৫০০) একজন গ্রীক দার্শনিক ও গণিজবিদ। বাংলা উইকিপিডিয়া

পিখাগোরাসের উপপাদ্য

A Must-read for contributors

পিখাগোরাস সবচেয়ে বেশি বিখ্যাভ পিখাগোরাসের উপপাদোর জন্য। গ্রীপ্টের জন্মের প্রায় ৬০০ বছর আগে পিখাগোরাস দর্শণ শাস্ত্রে প্রভুভ অবদান রাখেন। ভাকে বলা হব - সংখ্যার জনক। পিখাগোরাস প্রথম নিজেকে দার্শনিক হিসাবে দাবী করেন। প্লেটো, এরিপটল ও কোপারনিকাসের আনক অর্জন পিখ্যাগার্র্যাসর ধারণারই বিশ্রনি।

ভার কোনো দেখা পরে আর পাওয়া যায়নি ফলে ভার অবদান সম্পর্কে নিশ্চিত হওয়া যায়নি।

৫২৯ খ্রিস্টপর্বাদে তিনি ইতালীর ক্রোটোনা শহরে বসবাস শুরু করেন

পিখাগোরাসকে বলা যায় সংখ্যা দার্শনিক। তিনি বিশ্বাস করতেন পৃথিবীর সকল কিছর মল হলো সংখ্যা। সংখ্যার সঙ্গে নৈতিকতা, রঙ, শুভ-অশুভকে তিনি মিরিয়ে ফেলেন। তার বিশ্বাস হলো মানুষের আন্ন্যা অবিনশ্বর। ফলে হিন্দু, বৌদ্ধ ও অন্যান্য ক্ষেকটি ধর্মের মতো তিনিও জন্মান্তরবাদে বিশ্বাসী। পিখাগোরাস পথিবীর গোলকত্বে বিশ্বাস করতেন কারণ ভার ধারণা ছিল বত হলো পরম বক্ররেখা। সর্য, চন্দ্র বা গ্রহদের নিজেদের একটা গতি আছে।

[সম্পাদনা]



পিখাগোরাসের প্রভিক্তি

এই निवन्नि
ि अम्प्यूर्ण। आधिन हारेल ऽ
िएक मम्ब करां धारान।

প্লাক-সক্রেটীয় দার্শনিকবন্দ a - जा - x

সমকোণী ত্রিভ্তের অভিভ্তের ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফল অপর দুই বাহুর ওপর অঞ্চিত বর্গস্ক্রের ক্ষেত্রফলের সমষ্টির সমান।

Categories: Greek mathematician & philosopher

না • মেলিসোস

Born: 580 BC | Death: 500 BC

বিষয়দ্রেণীসমহ: অসম্পর্ণ । গ্লিক গণিতবিদ । গ্লিক দার্শনিক । খ্রিস্টপর্ব ৫৮০-এ জন্ম । খ্রিস্টপর্ব ৫০০-এ মত্য



#### **টে**ডিকিপিডিয়া वकि सक विश्वकाष

- পরিরমন প্রধান পাজা
- वाश्ना উইकिशिक्षिया সম্পদায
- সমসাম্বিক ঘটনা
- সাম্প্রতিক পরিবর্তনসমহ
- অজানা খেকোনো পর্তা
- সাহায্য ল দাল ককল

## অনুসন্ধান

চলো

অনসন্ধান

## হাজিয়ার

- সংযোগকারী পৃষ্ঠাসমৃহ সম্পর্কিত পরিবর্তন
- বিশেষ পৃষ্ঠাসমূহ
- সমকোণী ত্রিভুজের অভিভুজের ওপর অঞ্চিত বর্গক্ষেত্রের ক্ষেত্রফল অপর দুই বাহু ছাগার যোগ্য সংস্করণ স্থানী সংযোগ
- এই নিবন্ধটি উদ্ধৃত কক্ৰ

## অন্যান্য ভাষাসময়

- العربية =
- Asturianu Azərbaycan
- v. Žemaitėška
- Беларуская Български
- रेमात ठीत/विक्षिया মণিপরী
- Brezhonea Bosanski

। आगावना । राज्यानना पद्मन । राज्यान How to read and write in Unicode Bangla

A Must-read for contributors

#### <u> পিখাগোরাস</u> **Pythagoras**

শিশাগোরাস বা সামোসের শিশাগোরাস (প্রাচীন ত্রিক ভাষায় Ιωθαγόρας *পুখাগোরাস*) ( ত্রিস্টপূর্ব ৫৮০ - ত্রিস্টপূর্ব ৫০০) একজন গ্রীক দার্শনিক ও গণিজবিদ।

পিখাগোরাস সবচেয়ে বেশি বিখ্যাভ পিখাগোরাসের উপপাদোর জন্য। গ্রীপ্টের জন্মের প্রায় ৬০০ বছর আগে পিখাগোরাস দর্শণ শাস্ত্রে প্রভুভ অবদান রাখেন। ভাকে বলা হব - সংখ্যার জনক। পিখাগোরাস প্রথম নিজেকে দার্শনিক হিসাবে দাবী করেন। প্লেটো, এরিপটল ও কোপারনিকাসের আনক অর্জন পিখ্যাগার্র্যাসর ধারণারই বিশ্রনি।

ভার কোনো দেখা পরে আর পাওয়া যায়নি ফলে ভার অবদান সম্পর্কে নিশ্চিত হওয়া যায়নি।

৫২৯ খ্রিস্টপর্বাদে ভিনি ইভাদীর ক্রোটোনা শহরে বসবাস শুরু করেন।

পিখাগোরাসকে বলা যায় সংখ্যা দার্শনিক। তিনি বিশ্বাস করতেন পৃথিবীর সকল কিছর মল হলো সংখ্যা। সংখ্যার সঙ্গে নৈতিকতা, রঙ, শুভ-অশুভকে তিনি মিরিয়ে ফেলেন। তার বিশ্বাস হলো মানুষের আল্লা অবিনশ্বন। ফলে ছিল্ল বৌদ্ধ ও আন্তান কাষকটি ধার্মর মাজা তিনিও জন্মান্তরবাদে বিশ্বাসী। পিখাগোরাস পৃথিবীর গোলকত্বে বিশ্বাস করতেন কারণ তাঃ

Born on, PER Death on Cities of, LOC Countries of



পিখাগোরাসের প্রভিক্তি

এই निवन्निष्टि अमन्यूर्ण। आधिन हारेल अहिक ममृद्ध कराछ थातन।

a - जा - x

নিজেদের একটা গতি আছে।

পিখাগোরাসের উপপাদ্য

Categories: Greek mathematician & philosopher Born: 580 BC | Death: 500 BC

নো • মেপিসোস

বিষয়দ্রেণীসমহ: অসম্পর্ণ । গ্লিক গণিতবিদ । গ্লিক দার্শনিক । খ্রিস্টপর্ব ৫৮০-এ জন্ম । খ্রিস্টপর্ব ৫০০-এ মত্য

# Generating an Annotated List of Phrases (contd.)

# Steps

- Generate and annotate the title of each article
  - (i) Use category information to annotate the title
  - (ii) Use a small set of keywords to annotate the title

# Generating an Annotated List of Phrases (contd.)

NE Class	Keywords
PER	"born," "died," "one," "famous"
LOC	"city," "area," "population," "located," "part of"
ORG	"establish," "situate," "publish"

Table: Keywords for each named entity class

স্থানাংক: 22,5697, 88,3697

#### টেইকিপিডিয়া **बकिए सक विश्वकाय**

#### পরিভ্রমন

- প্রধান পাতা
- বাংলা উইকিপিডিয়া
- **म**म्लुपाय
- সমসাম্যিক ঘটনা
- সাম্প্রতিক পরিবর্তনসমহ
- অজানা (যকোনো পর্চা
- ॥ आहायर
- ল দাল কঞ্চল

#### অনুসন্ধান



#### হাজিযার

- সংযোগকারী পর্তাসমহ সম্পর্কিত পরিবর্তন
- বিশেষ পৃষ্ঠাসমূহ
- ছাপার যোগ্য সংস্করণ
- সামী সংযোগ
- এই নিবন্ধটি উদ্বৃত करूल

#### অন্যান্য ভাষাসমূহ

- Aragonés
- العربية و ্বসমীয়া
- Azərbaycan ৬ যাজায়াজ
- Беларуская ৭ জনপরিসংখ্যান
- Български । व्यक्तीन
- Brezhoneg

## How to read and write in Unicode Bangla A Must-read for contributors

সম্পাদনা কঞ্চন

## কলকাতা

উইকিপিডিয়া, মূক্ত বিশ্বকোৰ খেকে

নিবন্ধ আলোচনা

রাজধানী। হুগদী নদীর পর্ব জীরে অবস্থিত এই শহরের পৌর অঞ্চলর জনসংখ্যা ৫০ লক্ষের কিছ বেশি। ভবে কলকাতা ও ভার পার্রবর্তী জেলাগুলিতে বিস্তৃত কলকাতার মহানগরীয় অঞ্চলের জনসংখ্যা ১ কোটি ৪০ লক্ষের কাদাকাদি: এই জনসংখ্যার বিচারে কলকাতা ভারতের চতর্থ বহতম শহর ও ততীয় বহতম মেট্রোপলিটান বা

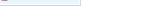
ইভিহাস

মহানগরীয অঞ্চল। ১৯১১ সাল পর্যন্ত কলকাতা রিটিশ ভারতের রাজধানী ছিল। একদা ভারতের আধুনিক শিক্ষা, বিজ্ঞান, শিব, সংস্কৃতি ও রাজনীতির কেন্দ্রভমি কলকাতা মহানগরে ১৯৫৪ সালের পর খেকে তীর রাজনৈতিক সংঘাত ও তার সঙ্গে সঙ্গে অর্থনৈতিক মন্দা দেখা দেখ। তার ২০০০ সালের পর থেকে এই শহর আবার অর্থিক ও বাণিজ্যিক সমদ্ধির পথে অগ্লসর হয় ও সাংস্কৃতিক হৃতগৌরব অনেকাংশে প্ররাধিকার করে। তবে ভারতের অন্যান্য মহানগরগুলির মতো লগরামূলজনিত দারিদ্রা, পরিবেশ দুষণ ও যানজাটের সমস্যা (থকে কলকাতাও একেবারে মূক্ত হতে পারেনি। কলকাতা শহরের প্রসিদ্ধি তার বৈপ্রবিক আন্দোলনগুলির জন্য। ভারতের স্বাধীনতা আন্দোলন এবং পরবর্তীতে বামপন্ধী বাণিজ্যিক ইউনিয়নগুলির আন্দোলন এ শহরের ইভিহাসের একটি বড অংশ। এর সঙ্গে সঙ্গে আধনিক

ভারতে কলকাতা *সাংষ্কৃতিক রাজধানী* ও *আনন্দ নগরী* (City of Joy *সিটি অব জয়*) নামে নন্দিত। রবীন্দ্রনাথ ঠাকর, রোনান্ড রস, সভাষ্টক্র বস, মাদার ভেরেসা, সভাজিৎ রায়, সভোক্তনাথ বস, স্থামী বিবেকানন্দ রাজা রামমোহন রায় সি ভি রামন-সহ বহু বিশ্ববিশ্রুত ব্যক্তিত্বের বাসভূমি এই নগরী ভার ঐতিহাসিক ঐতিহা ও সমৃদ্ধ সাংস্কৃতিক প্রেক্ষাপটের কারণে আজও বিশ্ববাসীর চোথে মর্যাদার আসনে অধিষ্ঠিত।



- ১ ইনিহাস ৩ ভূগোল
- ৩.১ শহারর গঠন ৪ অর্থনীতি
- েরাজধানী ও নগর প্রশাসন
- ৮ সংস্কৃতি ৯ শিক্ষা
- ১০ গণমাধ্যম



ক্ৰকাভা পশ্চিমবর • মাবন

क्रिकेविया सामावियान



**ব্যবাস:** 22.5697, 88.3697 IST (UTC+0:00)

আৰভন ১,৫০০ বৰ্গ কিমি (৫৭৯ sa mi) · উভৱা · ১ m (৩০ ft)

ক্ৰকাভা (অনা

कनभर शा 8,६४०,६88 (२००५) েদটো ১,১২ন্ট (এক্সপ্রেশন ক্রটি: ধ-এর ফেন্য অপজ্লবন্ড নেই। 🤉 😁



**बकिए सक विश्वकाय** 

How to read and write in Unicode Bangla

সম্পাদনা কঞ্চন

Must-read for contributors

কলকাতা Kolkata

নিবন্ধ আলোচনা

উথকদিভিয়া, মূক্ত বিশ্বকোষ খেকে

ষ্যালাংক: 22.5697, 88.3697

#### পরিভ্রমন

- প্রধান পালা
- वाश्ना उँशैकित्रिछिया
- সম্প্রভিক পরিবর্তনসমহ
- মাণ্লাভক পারবভনসমূ
   অজানা মেকোনো পর্চা
- आहास
- » দান কঞ্চন
- m allet delta

## অনুসন্ধান



## **চলো** অনুসন্ধান হাতিযার

#### NIO.NA

- সংযোগকারী পৃষ্ঠাসমূহ
   সম্পর্কিত পরিবর্তন
- বিশেষ পৃষ্ঠাসমূহ
- ছাপার (যাগ্য সংস্করণ
- ব্যায়ী সংযোগ

   এই বিবলটী উচ্চত
- এই নিবন্ধটি উদ্ভ করুন

#### অন্যান্য ভাষাসমূহ

- Aragonés العربية
- = الحربية = = अप्रश्रीया (\*
- Azərbaycan ৬ যাভায়াভ
- Беларуская । ৭ জনপরিস
- н Былгарски н бебе
- Brezhoneg

কৰকাভা (ইংরেজি ভাষান: Kolkata, পূর্বে • Calcutta <sup>আঞ্জনকা</sup>) ভারতের পূর্বাঞ্গীন রাজা পশ্চিনবাসের রাজমানী। হুগদী নদীর পূর্ব জীরে অবস্থিত এই শহরের পৌর অঞ্চলের জনসংগা ৫০ লাক্তর কিছু (বিশি ভাবে কলকাভা ও ভার পার্ববাজী লোগাধিটিভ বিস্কৃত কলকাভার নহানগরীন অঞ্চলের অনসংগা ১ (কাটি ৪০ শাহরর কাছাকাছি, এই অলসংখ্যার বিচারে কলকাভা ভারতের চকুপ বৃহত্তম শহরে ও ভৃতীর বৃহত্তম দেটোপদিচান বা মহামান্ত্রীয় অঞ্চল।

ইভিহাস

১৯১১ সাল পর্যন্ত কলকাতা ব্লিটিশ ভারভের রাজখানী ছিল। একদা ভারভের আধুনিক শিরা, বিজ্ঞান, শিব, সংস্কৃতি ও রাজনীভির কেন্দ্রভূমি কলকাতা মহানগরে ১৯৫৪ সালের পর থেকে জীর রাজনৈভিক সংঘাভ ও ভার সঙ্গে সঙ্গে অর্থনৈভিক মলা দেখা দেখা ভাব ২০০০ সালের পর থেকে এই শহর অব্যার আর্থিক ও বালিজিকে সমৃদ্ধির পৃথ অগ্রসর হব ও সাংস্কৃতিক কভাগীরব আবেকাংশ পুররাধিকার করে। ভাব ভারভের অলান্দ্র মহালার্ডন মাতা নকারান্ত্রনাভিক কলাতা পার্যন্ত্র প্রতিবিশ্ব শৃবণ ও বালজার্ত্তর মহাসা খেকে কলকাতা ব্যাহনিত আবেলান্দ্র হতে পার্রেন। কলকাতা শার্রের প্রমিট্ড ভার থেমবিক আন্দোলন করা। তারভের প্রাধীকতা আব্দোলন এবং পরবর্তীতে বাসপন্থী বালিজাকে ইউনিয়নভালির আন্দোলন এ শহরের ইভিয়াসের একটি বড় অংশ। এর সঙ্গে সঙ্গে আধুনিক

ভারতে কবকাতা সাংস্কৃতিক রাজখানী ও অনন্দ নগারী (City of Joy সিটি অব করে) নামে নদিত। রবীচনাখ ঠাকুর, হোলাভ রম, দুড়াযথত রমু মাদার তেরেমা, মতাতিম রাম, মাতাচনাম বামু, যানী বিবেচনাশ রাজা রামামানের রাম বি চামান-ম যব বিবিহিত্ত বাহিত্রত মাত্রাক্তর মাত্রান্ত এই করি ভার ঔতিহাসিক ঐতিহা ও সমূদ্দ সাংস্কৃতিক প্রস্কাণ্টের কারণে আকও বিহ্ববাসীর চোখে মর্মানার আদান অমিষ্টিভ।

### সূ**টিপত্র** [আড়ালে রাখো] ১ নামকরণ

- ২ ইভিহাস
- ৬ ভূগোল ৬.১ শহরের গঠন
- ৪ অর্থনীতি
- েরাজধানী ও নগর প্রশাসন
- ৭ জনপরিসংখ্যান
- ৮ সংস্কৃতি ১ শিক্ষা
- ১০ গণমাধ্যম





ক্ৰকাভা

পশ্চিমবর • মাবন

क्रिकेविया सामावियान

মুর IST (UTC+৫:৩০)

আমেতন ১,৫০০ বৰ্গ কিমি (৫৭৯ sq mi) • উচ্চতা • ৯ m (৩০ ft)

(बना कनकां (बना

**জনসংখ্যা** ৪,৫৮০,৫৪৪ (২০০১)

্ষনত্ব ...

ব্যাহার বিষয়ে বিষয় বি







# Generating an Annotated List of Phrases (contd.)

## Steps

- Generate and annotate the title of each article
  - (i) Use category information to annotate the title
  - (ii) Use a small set of keywords to annotate the title
- ② Getting more location names
- Generating and annotating the tokens in the titles

# Generating an Annotated List of Phrases (contd.)

# Annotating the tokens in the titles

- Assign each token the same NE label as that of its title
- Ambiguous case:

Solution: Label "Anna" with its most frequent NE class

## Example

Smith College is in Massachusetts

## Example

Smith College is in Massachusetts

## Example

Smith College is in Massachusetts

ORG ORG



# Example Smith College is in Massachusetts ORG ORG OTHERS OTHERS LOC

# Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- Pipelined NER with New Features
  - Evaluation
- 4 Joint Model for POS Tagging and NEF
  - Evaluation

# POS Tagging

## Goal

Investigate whether the two new features can improve a baseline supervised POS tagger

## Experimental setup

- Corpus: 78K word tokens
- Tagset: IIIT Hyderabad's POS tagset with 26 tags
- Learning algorithm: CRF
- 5-fold cross-validation (CV)

# **Features**

## Baseline

word n-grams, pseudo-affixes

# Results

Experiment	Overall	Seen	Unseen
Baseline	89.8	92.9	72.0

Table: 5-fold cross-validation accuracies

# **Features**

## Baseline+Induced Affixes

- word n-grams, pseudo-affixes
- induced affix n-grams

# Results

Experiment	Overall	Seen	Unseen
Baseline	89.8	92.9	72.0
Baseline+Induced Affixes	90.5	93.3	74.6

Table: 5-fold cross-validation accuracies

# **Features**

## Baseline+Induced Affixes+Wiki

- word n-grams, pseudo-affixes
- induced affix n-grams
- induced Wiki n-grams

Experiment	Overall	Seen	Unseen
Baseline	89.8	92.9	72.0
Baseline+Induced Affixes	90.5	93.3	74.6
Baseline+Induced Affixes+Wiki	90.8	93.5	75.5

Table: 5-fold cross-validation accuracies

## Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- Pipelined NER with New Features
  - Evaluation
- 4 Joint Model for POS Tagging and NEF
  - Evaluation

#### **NER**

#### Goal

Investigate whether the two features can improve a pipelined NE recognizer that uses induced POS tags as features

#### Experimental setup

- Corpus: 78K word tokens
- Learning algorithm: CRF
- IOB convention
- 5-fold cross-validation (CV)

## **Features**

#### Baseline

word n-grams, pseudo-affixes, POS n-grams

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Person	66.1	74.0	69.9
Organization	29.8	44.9	35.8
Location	52.6	80.4	63.6

#### **Features**

#### Baseline+Induced Affixes

- word n-grams, pseudo-affixes, POS n-grams
- induced affix n-grams

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Person	66.1	74.0	69.9
Organization	29.8	44.9	35.8
Location	52.6	80.4	63.6
Baseline+Induced Affixes	60.4	73.3	66.2
Person	65.7	72.6	69.0
Organization	31.7	46.4	37.7
Location	51.4	80.0	62.6

### **Features**

#### Baseline+Induced Affixes+Wiki

- word n-grams, pseudo-affixes, POS n-grams
- induced affix n-grams
- induced Wiki n-grams

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Person	66.1	74.0	69.9
Organization	29.8	44.9	35.8
Location	52.6	80.4	63.6
Baseline+Induced Affixes	60.4	73.3	66.2
Person	65.7	72.6	69.0
Organization	31.7	46.4	37.7
Location	51.4	80.0	62.6
Baseline+Induced Affixes+Wiki	63.2	75.1	68.7
Person	66.4	75.1	70.5
Organization	30.7	43.8	36.1
Location	60.0	79.6	68.5

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Person	66.1	74.0	69.9
Organization	29.8	44.9	35.8
Location	52.6	80.4	63.6
Baseline+Induced Affixes	60.4	73.3	66.2
Person	65.7	72.6	69.0
Organization	31.7	46.4	37.7
Location	51.4	80.0	62.6
Baseline+Induced Affixes+Wiki	63.2	75.1	68.7
Person	66.4	75.1	70.5
Organization	30.7	43.8	36.1
Location	60.0	79.6	68.5

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Person	66.1	74.0	69.9
Organization	29.8	44.9	35.8
Location	52.6	80.4	63.6
Baseline+Induced Affixes	60.4	73.3	66.2
Person	65.7	72.6	69.0
Organization	31.7	46.4	37.7
Location	51.4	80.0	62.6
Baseline+Induced Affixes+Wiki	63.2	75.1	68.7
Person	66.4	75.1	70.5
Organization	30.7	43.8	36.1
Location	60.0	79.6	68.5

## Outline

- Two New Linguistic Features
  - Induced Affixes
  - Semantic Classes from Wikipedia
- POS Tagging with New Features
  - Evaluation
- Pipelined NER with New Features
  - Evaluation
- Joint Model for POS Tagging and NER
  - Evaluation

## Joint Model

#### Goal

Jointly predict a word's POS and NE tag to help avoid error propagation

#### Model assumption

A Bengali word is part of an NE if and only if it is a proper noun. For our evaluation corpus, this assumption is correct 97.3% of the time.

# Joint Model (contd.)

#### Model

- Jointly predicts a word's POS and NE tag
- Trained using CRF
  - Features: POS and NE taggers' features minus NE tagger's POS-related features
  - Joint tag: If a word is not a proper noun, its class is simply its POS tag. Otherwise, its class is its NE tag

#### Example:

Dhaka is the capital of Bangladesh

B-LOC VBZ DT NN IN B-LOC

# **NER Results for Joint Modeling**

Experiment	R	Р	F
Baseline	54.7	81.7	65.5
Baseline+Induced Affixes	56.7	88.9	69.3
Baseline+Induced Affixes+Wiki	61.7	86.3	71.9

Table: 5-fold cross-validation joint modeling results for NER

# Comparison between Joint and Pipelined Model

Experiment	R	Р	F
Baseline	54.7	81.7	65.5
Baseline+Induced Affixes	56.7	88.9	69.3
Baseline+Induced Affixes+Wiki	61.7	86.3	71.9

Table: 5-fold cross-validation joint modeling results for NER

Experiment	R	Р	F
Baseline	60.9	74.4	67.0
Baseline+Induced Affixes	60.4	73.3	66.2
Baseline+Induced Affixes+Wiki	63.2	75.1	68.7

# Summary

- Investigate the pipelined NER architecture and proposed two new features:
  - induced affixes
  - semantic class information from Wikipedia
- Proposed a joint model for learning POS tagging and NER simultaneously