# Persian Language, is Stemming Efficient?

Ljiljana Dolamic, Jacques Savoy

Computer Science Department

University of Neuchatel, Switzerland
www.unine.ch/info/clef/

Université
de Neuchâtel **unine**

# Outline

- Introduction
- Persian Language
- Indexing strategies
- IR models
- Evaluation
- Conclusion

# Introduction

- Persian language -50 million native speakers
- Test-collection
  - University of Teheran, CLEF 2008
  - Hamshahri (1996-2002)
  - 166,477 documents
  - 100 topics
    - 501-550 – training topics for CLEF 2008
    - 551-600 – test topics for CLEF 2008
  - Relevance judgement

# Introduction

☐ **&lt;topic lang="fa"&gt;**

&lt;identifier&gt;**10.2452/584-AH**&lt;/identifier&gt;

&lt;title&gt;**تحولات قیمت نفت**&lt;/title&gt;

&lt;description&gt; **چگونه جهان و ایران بازار در نفت قیمت تحولات**
**است؟**&lt;/description&gt;

&lt;narrative&gt; **روي آن تاثیر و جهاني بازارهاي در نفت قیمت تحولات آمار خواهم مي**
**چیست؟ تورم و اقتصاد**&lt;/narrative&gt;

**&lt;/topic&gt;**

# Outline

# Persian Language

- ☐ Indo-European language
  - ■ Distinctly related to European languages
- ☐ Written in Arabic script (28 letters)
  - ■ Additional four letters
    - ☐ گ ژ چ پ
      - ■ پدر – father (pedar)
      - ■ بچّه – child (bačče)
      - ■ ژاله – dew (žāle)
      - ■ مرگ – death (marg)

# Persian language – writing system

- Written from right to left
- Four different forms to a letter
  - ت – teh
    - توت – <u>unconnected</u> and <u>initial</u>
    - رست – <u>final</u>
    - چتر – <u>medial</u>
- Letters ا, د, ذ, ر, ز, ژ and و do not unite with the letter following

# Persian language - vowels

- ☐ All letters are consonants
  - ■ ا و ی – used to express vowel sounds 'دو' (two)
- ☐ Short vowels
  - ■ Vowel signs 'دَر', 'پُشت', 'دِل'
- ☐ Long vowels
  - ■ Vowel sign + ا و ی - 'کَار'
- ☐ Other signs
  - ■ ـَّ , ـْ , ـٔ

# Persian language

- ☐ Definite article
  - ■ كتاب - "book" or "the book"
  - ■ آب بيار (bring water) - آب را بيار (bring the water)
  - ■ مردى كه - the man who
  - ■ پسره - son (in question)

# Persian Language

□ Indefinite article

■ كتاب (book) – كتابى (a book)

■ يک كتاب –a book

■ يک پادشاهى – a king

■ خانهٔى – a house

■ نامهٔ –a letter

■ مردهائى – some men

# Persian Language

- ☐ No grammatical gender
  - ■ مرد – man, زن – woman
  - ■ شیر نر – lion, شیر ماده – lioness
  - ■ Arabic words
    - ☐ Danser: رقاص (male) - رقاصه (female)

# Persian Language

- ☐ Declination
  - ■ Accusative: زن را (the woman), مرد را (the man)
  - ■ Genitive: خانهٔ مرد (the man's house)

    پسر مرد (the man's son)

- ☐ Prepositions for expressing relations
  - ■ از مرد (from the man)

# Persian Language

- ☐ Plural endings:
  - ■ ان (animate), ها (inanimate)
    - ☐ بندگان ← بنده, آقایان, پدران
    - ☐ مردها, خانه‌ها, گلها
  - ■ Arabic plurals:
    - ☐ Sound plurals: ین (masculine) , ات (feminine)
    - ☐ Broken plurals: قلوب ← قلب
- ☐ Adjectives
  - ■ Substantives: plural endings بزرگان (the great)

# Persian language - problems

- Vowel signs omitting
  - man – (مَرد) مرد – dead vs. (مُرده) مرده
- ZWNJ
  - خانه‌ها or خانه ها and not خانه‌ها
- Arabic vs. Persian letters
  - ي (\u064A), ى (\u0649) vs. ی (\06CC)
  - یـ ,ی ,ی vs. یـ ,یـ ,ي ,ي

# Indexing strategies - stemming

- "none"
- Light stemmer - www.unine.ch/info/clef/
- perstem
  - Jon Dehdari
  - http://sourceforge.net/projects/perstem/
  - *"If these regular expressions are readable to you, you need to check-in to a psychiatric ward!"*
- 5-gram

# Indexing strategies - stoplist

- Available at : www.unine.ch/info/clef/
  - 881 term
  - Non content baring terms
  - Suffixes and prefixes
    - Ex.: plural suffixes ان or ها

# Outline

- Introduction

- Persian Language

- Indexing strategies

- **IR models**

- Evaluation

- Conclusion

# IR models

- Vector-space
  - *Lnu-ltc*
  - *tf-idf*
- Probabilistic
  - Okapi
  - DFR-*PL2*
  - Language model LM

# Outline

- Introduction

- Persian Language

- Indexing strategies

- IR models

- Evaluation

- Conclusion

# Evaluation

| T query | none | light | perstem | 5-gram |
|---|---|---|---|---|
| *tf idf* | 0.2966 | 0.2625* | 0.2799* | 0.2814 |
| *Lnu-ltc* | 0.4533 | 0.4277* | 0.4369 | 0.4107* |
| Okapi | 0.4811 | 0.4535* | 0.4610* | 0.4511 |
| DFR-PL2 | **0.4939** | **0.4693*** | **0.4750*** | **0.4423*** |
| LM | 0.4348 | 0.4000* | 0.4113* | 0.3892* |
| over "none" | | -6.79% | -4.43% | -8.57% |

# Evaluation

- Topic #525 "Seasonal Diseases"
  - فصلي (seasonal) → فصل (season, term, article )
  - فصل كم كاري - layoff
- Topic #522 "Teheran metro project"
  - متروي (metro)
  - مِثْروي
  - مِتر (meter)

# Evaluation

| T query | none | | perstem | |
|---|---|---|---|---|
| | stoplist | no stoplist | stoplist | no stoplist |
| *tf idf* | 0.2966 | 0.2449* | 0.2799 | 0.2341* |
| *Lnu-ltc* | 0.4533 | 0.4519 | 0.4369 | 0.4363 |
| Okapi | 0.4811 | 0.4806 | 0.4610 | 0.4591 |
| DFR-PL2 | **0.4939** | **0.4888** | **0.4750** | **0.4747** |
| LM | 0.4348 | 0.4186* | 0.4113 | 0.4034 |
| over "none" | | 2.74% | | 3.47% |

# Evaluation

- □ Avdl
  - ■ 2.79(stop list) vs. 3.05 (no stop list)
- □ TD queries
  - ■ 5.94 (stop list) vs. 9.62 (no stop list)
  - ■ TD queries, light stemmer

| TD query | *tf idf* | Okapi | DFR-PL2 | Over no stop list |
|----------|----------|-------|---------|-------------------|
| *stoplist* | 0.2744 | 0.4559 | 0.4785 | 11.55% |
| *no stoplist* | 0.2176 | 0.4110 | 0.4576 | |

# Outline

- Introduction

- Persian Language

- Indexing strategies

- IR models

- Evaluation

- **Conclusion**

# Conclusion

- ☐ Main aspects of Persian morphology
- ☐ Indexing strategies
  - ■ Stemming hurts retrieval effectiveness
  - ■ Effect of the stop list removal depends on query size
- ☐ IR models
  - ■ *Divergence from randomness*

# Persian Language,
# is Stemming Efficient?

## THANK YOU!

Ljiljana Dolamic, Jacques Savoy

Computer Science Department

University of Neuchatel, Switzerland
www.unine.ch/info/clef/

Université
de Neuchâtel **uni**ne