

# LINGÜÍSTICA COMPUTACIONAL

## TOKENIZING

Antonio Frías Delgado

Área de Lógica y Filosofía de la Ciencia  
Departamento de Historia, Geografía y Filosofía  
Universidad de Cádiz

# SEGMENTACIÓN: BIBLIOGRAFÍA

# SEGMENTACIÓN: BIBLIOGRAFÍA

- What is a word, What is a sentence? Problems of Tokenization  
Gregory Grefenstette, Pasi Tapanainen

# SEGMENTACIÓN: BIBLIOGRAFÍA

- What is a word, What is a sentence? Problems of Tokenization  
Gregory Grefenstette, Pasi Tapanainen
- Text preprocessing  
David Palmer

# SEGMENTACIÓN: BIBLIOGRAFÍA

- What is a word, What is a sentence? Problems of Tokenization  
Gregory Grefenstette, Pasi Tapanainen
- Text preprocessing  
David Palmer
- Tokenisation and sentence segmentation  
David Palmer

# SEGMENTACIÓN: RECURSOS

# SEGMENTACIÓN: RECURSOS

NLTK

# SEGMENTACIÓN: RECURSOS

## NLTK

- `import nltk`



# SEGMENTACIÓN: RECURSOS

## NLTK

- `import nltk`
- `nltk.word_tokenize(texto, language='spanish')`

# SEGMENTACIÓN: RECURSOS

## NLTK

- `import nltk`
- `nltk.word_tokenize(texto, language='spanish')`
- `nltk.wordpunct_tokenize(texto)`

# SEGMENTACIÓN: RECURSOS

## NLTK

- `import nltk`
- `nltk.word_tokenize(texto, language='spanish')`
- `nltk.wordpunct_tokenize(texto)`
- `nltk.sent_tokenize(texto, language='spanish')`

# NLTK: BROWN CORPUS

- `import nltk`

# NLTK: BROWN CORPUS

```
■ import nltk  
■ from nltk.corpus import brown
```

# NLTK: BROWN CORPUS

```
■ import nltk  
■ from nltk.corpus import brown  
■ text_brown=brown.raw()
```

# NLTK: BROWN CORPUS

```
■ import nltk
■ from nltk.corpus import brown
■ text_brown=brown.raw()
■ oraciones_brown=brown.sents()
```

# NLTK: BROWN CORPUS

```
■ import nltk
■ from nltk.corpus import brown
■ text_brown=brown.raw()
■ oraciones_brown=brown.sents()
■ palabras_brown=brown.words()
```