# **ALTAMURA SANTE - Pos Tagger per lingue morte: Latino e Greco**

## 0. Introduzione

L'obiettivo dell'esercitazione è quello di implementare un Pos Tagger statistico per lingue morte (Latino e Greco) basato su **HMM** (Hidden Markov Model), il quale molto spesso viene utilizzato nella analisi delle sequenze.

L'operazione di Pos Tagging consiste in:

- Input: una frase, vista come una sequenza di parole
- **Output**: una lista di Pos Tag. Ogni tag è associato ad una parola della frase presa in input dall'algoritmo.

Questo problema, secondo lo standard del machine learning, è stato affrontato in 3 fasi, ma la definizione del modello (Modelling) è già nota:

• **Modelling**: definizione del **modello matematico** del problema L'obiettivo è trovare la sequenza di tag  $\hat{t}_1^n$  data la sequenza di parole osservabili  $w_1^n$ , che massimizza la distribuzione di probabilità P:

$$\hat{t}_1^n = rgmax_{t_1^n} P(t_1^n, w_1^n)$$

Per rendere operazionale questo calcolo è possibile sfruttare la regola di Bayes in modo tale da ottenere una equazione che avrà più probabilità da calcolare, ma **approssimabili**. In questo modo si ottiene una nuova equazione che approssima la precedente:

$$\hat{t}_1^n = rgmax_{t_1^n} P(t_1^n, w_1^n) pprox \prod_{i=1}^n P(w_i \mid t_i) P(t_i \mid t_{i-1})$$

Distinguiamo le due probabilità:

$$P(w_i \mid t_i) \in P(t_i \mid t_{i-1})$$

le quali sono rispettivamente le **probabilità di transizione** ed **emissione** calcolate nella fase di Learning, che verranno trattate nel capitolo successivo insieme agli algortimi e strutture dati impiegati rispettivamente per ottenerli e memorizzarli

- **Learning**: apprendere i parametri da un dato corpus, che verranno utilizzati nella fase di Decoding
- Decoding: implementazione dell'algoritmo di Pos Tagging

# 1. Implementazione del Learning

Le probabilità di transizione ed emissione vengono calcolate indipendentemente mediante i **train set** del greco e latino.

## 1.1 Probabilità di transizione

La probabilità di transizione rappresenta la probabilità di un tag dato il tag precedente. I **conteggi** vengono definiti rispettivamente dalla formula della probabilità:

$$P(t_i,t_{i-1}) = rac{C(t_{i-1},t_i)}{C(t_{i-1})}$$

 $C(t_{i-1}, t_i)$  = numero di volte che il tag  $t_{i-1}$  compare prima del tag  $t_i$  all'interno del train set

```
# sent_id = train-s1749
# text = quia inter nos taliter convinet a Rachisindo notario iscrivere rogavit.
# reference = document_id='73:20'-span='110'
1  quia  quia  SCONJ
2  inter  inter  ADP i-1
3  nos  nos  PRON i
4  taliter  taliter  ADV
5  convinet  conuenio  VERB
```

Ogni parola della frase ha un attributo *upos* e un attributo *form* mediante i quali è possibile ricavare rispettivamente il Pos Tag e la forma (token) della parola stessa.

#### 1.2 Matrice di transizione

Per la memorizzazione delle probabilità di transizione calcolate come descritto nel capitolo precedente è stata utilizzata una **matrice** etichettata sia sulle righe che sulle colonne con i Pos Tag per poter accedere alle probabilità in maniera semplice, specificando i due tag relativi ad essa.

Le **matrici di transizione** cambiano in base al linguaggio utilizzato per il training perchè l'insieme dei Pos Tag relativi al train set del latino è diverso da quello relativo al greco antico.

Infatti avremo una lista di tag possibili (*possible\_tags*) a cui è stato aggiunto lo stato iniziale 'START' e lo stato finale 'END'. Quest ultimi NON sono associati in generale a nessuna osservazione (parola) ma saranno utili nella fase di Learning per il calcolo delle probabilità di transizione iniziali e finali utilizzate successivamente nella fase di Decoding.

```
Pos Tags Greco: ['START','ADJ','ADP', 'ADV', 'CCONJ', 'DET', 'INTJ', 'NOUN', 'NUM', 'PART', 'PRON','SCONJ', 'VERB', 'X', 'PUNCT','END']
```

Di seguito l'algoritmo compute\_transition\_matrix nel file **learning2.py** utilizzato per popolare la matrice di transizione. Il calcolo viene effettuato per ogni coppia di tag presente nella lista **possible\_tags**:

```
sentence_n = 0
   for sentence in parse_incr(train):
       sentence n += 1
       for i in range(len(sentence)):
           word_before = sentence[i-1]
           word = sentence[i]
           if i == 0:
               if word["upos"] in count_initial_dict.keys():
                   count_initial_dict[word["upos"]] =
count_initial_dict[word["upos"]] + 1
           if (word_before["upos"], word["upos"]) in
transition_counter_dict.keys() and i != 0:
               transition_counter_dict[(word_before["upos"], word["upos"])] =
transition_counter_dict[(word_before["upos"], word["upos"])] + 1
           if word["upos"] in counter_dict.keys():
               counter_dict[word["upos"]] = counter_dict[word["upos"]] + 1
           if i == len(sentence) - 1:
               if (word["upos"], 'END') in transition_counter_dict.keys():
                   transition_counter_dict[(word["upos"], 'END')] =
transition_counter_dict[(word_before["upos"], word["upos"])] + 1
   #-----#
   #probabilità di transizione iniziali
   for i,t in enumerate(possible_tags):
       transition_matrix[0][i] = count_initial_dict[t]/sentence_n
   #probabilità di transizione intermedie
   for i,t1 in enumerate(possible_tags):
       for j,t2 in enumerate(possible_tags):
           if i \ge 1 and j \ge 1 and i < (len(possible_tags) - 1):
                transition_matrix[i][j] =
transition_counter_dict[(t1, t2)]/counter_dict[t1]
   train.seek(0)
   return transition_matrix
```

- Nella **prima fase** vengono inizializzati i tre dizionari utilizzati per memorizzare i conteggi relativi ai tag. In particolare abbiamo:
  - counter\_dict: dizionario relativo ai conteggi dei tag singoli all'interno del train set. Sarà formato da coppie (TAG, conteggio) tale che il valore 'conteggio' indica quante volte il TAG compare all'interno del train set.
  - count\_initial\_dict: dizionario relativo ai conteggi dei tag singoli che compaiono per primi nelle sentence del train set. Sarà formato da coppie (TAG, conteggio) tale che il valore 'conteggio' indica quante volte il TAG è associato alla prima parola delle sentence all'interno del train set.
  - count\_initial\_dict: dizionario relativo ai conteggio delle coppie di tag che compaiono uno dopo l'altro nelle sentence del train set. Sarà formato da coppie ((TAG1, TAG2), conteggio), tale che il valore 'conteggio' indica quante volte il TAG1 compare immediatamente prima del TAG2 nelle sentence del train set.

In precedenza viene inizializzata la matrice di transizione come una matrice di zeri, quindi vuota

• Nella **seconda fase** i tre dizionari vengono popolati con i relativi conteggi. Viene analizzata ogni sentence del train set in modo tale da effettuare tutti i conteggi e memorizzarli all'interno dei dizionari.

Nella terza fase la matrice di transizione viene riempita con le probabilità di transizione. <u>La probabilità del tag2 dato il tag1 è memorizzata nella cella delle matrice avente come riga il tag1 e come colonna il tag2.</u>

In particolare vengono effettuate due operazioni:

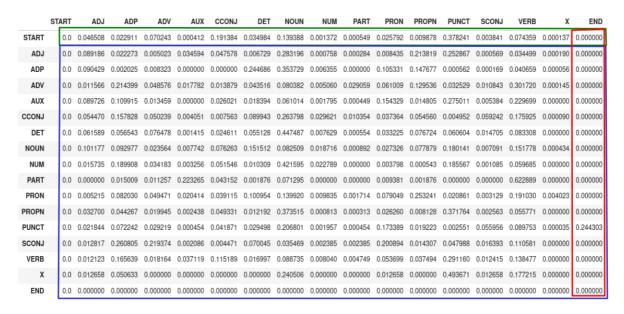
- 1. Viene popolata la prima riga della matrice (**riquadro verde**). di transizione con le probabilità di transizione iniziali. Ogni cella è calcolata dividendo il numero di volte che il tag relativo alla colonna è associato alla prima parola di ogni sentence per il numero di sentence totali.
- 2. Vengono popolate le righe successive alla prima (**riquadro blu**) con le probabilità di transizione calcolate come è stato descritto nel capitolo precedente, dividendo il numero di volte che il tag1(riga) compare immediatamente prima del tag2 (colonna) per il numero di volte che compare il tag1.

Dunque, ogni riga della matrice di transizione (tranne l'ultima) rappresenta una distribuzione di probabilità e la **somma dei suoi valori è sempre pari a 1**.

Nell'ultima colonna della matrice di transizione sono memorizzate le probabilità di transizione finali **(riquadro rosso)**. E' importante notare che sono tutte pari a 0 tranne quella relativa a PUNCT; questo accade perchè tutte le sentence del train set terminano con un simbolo di punteggiatura.

Successivamente la matrice viene trasformata in un **DataFrame** per poter etichettare le righe e le colonne con i Pos Tag.

Un esempio di matrice di transizione relativa al train set per il Latino:



#### 1.3 Probabilità di emissione

La probabilità di emissione rappresenta quanto un certo tag sia una certa parola. Anche in questo caso i conteggi vengono definiti dalla formula per il calcolo della probabilità:

$$P(w_i,t_i) = rac{C(t_i,w_i)}{C(t_i)}$$

 $C(t_i, w_i)$  = numero di volte che la parola  $w_i$  è taggata con il tag  $t_i$  nel train set

 $C(t_i)$  = numero di volte che il tag  $t_i$  compare nel train set

## 1.4 Dizionario di emissione

Per la memorizzazione delle probabilità di emissione è stato utilizzato un **dizionario** in quanto una matrice di emissione avrebbe avuto tante colonne quante sono le parole nel train set e non sarebbe una struttura dati ottimale.

La funzione compute\_emission\_probabilities nel file **learning2.py** restituisce tre dizionari:

- il dizionario di emissione, descritto a breve.
- il **dizionario delle parole** con il relativo conteggio all'interno del train set. Servirà nella fase di Decoding per verificare se una parola della sentence analizzata è conosciuta o no
- il dizionario delle parole conteggiate con i pos tag. E' utile nell'algoritmo di Baseline.

La funzione prende in input il train set (del Latino o del Greco) e restituisce come primo elemento **emission\_dict**, un dizionario che avrà come chiavi delle coppie **(word,tag)** e come valori le **probabilità di emissione** relative alle coppie.

Di seguito, un piccolo estratto del dizionario di emissione associato al train set del Latino:

```
{...,('Dei', 'PROPN'): 0.057208953357509064, ('nomine', 'NOUN'):
0.01572524239062274, ('regnante', 'VERB'): 0.015205799033494418, ('domno',
'NOUN'): 0.014205778785393855, ('nostro', 'DET'): 0.035255029840644804,
('Carulo', 'PROPN'): 0.004689258471926972, ('rege', 'NOUN'):
0.0035695335487916646, ('Francorum', 'NOUN'): 0.0027977425112150887,...}
```

# 2. Implementazione del Decoding

La fase di Decoding permette di capire qual è l'algoritmo che permette di applicare le probabilità apprese nella fase di Learning per poter restituire la sequenza di tag ottimale per una determinata frase in input.

La soluzione è quella di applicare la **dynamic programming** e passare da una complessità esponenziale ad una **polinomiale** attraverso l'**approssimazione Markoviana**.

Le due assunzioni fondamentali dell'HMM sono:

- <u>La probabilità che una parola appaia dipende solo dal suo stesso tag ed è indipendente dalle parole vicine e dagli altri tag.</u>
- <u>La probabilità di un tag dipende solo dal tag precedente anzichè che dall'intera sequenza di tag precedenti</u>.

# 2.1 Algoritmo di Viterbi

L'algoritmo di Viterbi implementato è una **leggera "variante"** di quello tradizionale in quanto utilizza diverse strutture dati, ma con la stessa finalità. <u>L'algoritmo, con le sue funzioni di supporto, è implementato nel file **viterbi.py**.</u>

Input dell'algoritmo:

- sentence\_tokens : una lista di parole della frase analizzata
- possible\_tags: lista dei possibili tag relativi dal train set utilizzato (Greco o Latino)
- transition\_matrix: matrice di transizione che memorizza tutte le probabilità di transizione
- emission\_probabilities: dizionario delle probabilità di emissione
- count\_word: dizionario delle parole con i relativi conteggi nel train set
- smoothing\_strategy: strategia di smoothing utilizzata, relativa al train set
- oneshot\_words\_tag\_distribution: distribuzione probabilistica dei tag relativi alle parole che compaiono una sola volta nel dev set

#### Strutture dati principali:

```
viterbi_matrix = np.zeros((len(possible_tags), len(sentence_tokens)))
backpointer = dict()
```

• viterbi\_matrix: matrice di dimensione [possible\_tags] \* |sentence\_tokens], che sono rispettivamente il numero di tag possibili dai quali sono stati rimossi 'START' ed 'END' ed il numero di termini della frase in input.

Una differenza importante con l'algoritmo standard di Viterbi visto a lezione è la struttura della matrice di Viterbi: non è stato considerato necessario ampliare la dimensione della matrice con le righe corrispondenti ai tag di 'START' ed 'END' per due motivi:

- 1. gli stati corrispondenti a questi tag non vengono calcolati. Ciò che è importante, invece, sono le probabilità presenti nella matrice di transizione.
- 2. viene restituita dall'algoritmo solo la sequenza più probabile di tag, non la sua probabilità, la quale non viene memorizzata in nessuno stato finale.

Ogni cella della matrice che chiameremo 'stato' ha l'utilità di memorizzare la probabilità più alta di una sottosequenza della sequenza di tag, questo ci permette di ridurre la complessità perchè NON vengono calcolate iterativamente tutte le possibili sequenze per ottenere la soluzione migliore sfociando in un'esplosione combinatoria.

• backpointer: dizionario di dizionari. Ad ogni chiave (colonna) corrisponde un dizionario che rappresenta una colonna. Ogni colonna è rappresentata da un insieme di chiavi (righe) a cui corrisponde il puntatore alla riga della colonna precedente.

Questa struttura dati è utile per memorizzare il riferimento di ogni stato allo stato precedente che ha dato maggior contributo nel calcolo del suo valore di Viterbi.

L'algoritmo è diviso in 4 fasi:

#### 1.Inizializzazione della prima colonna

Viene inizializzata la prima colonna della matrice di Viterbi: per ogni stato rappresentato da un tag viene calcolata la somma dei logaritmi tra la probabilità di transizione iniziale relativa al tag e la probabilità di emissione relativa alla coppia (parola, tag). Se la probabilità di transizione o quella di emissione è pari a 0, la si transforma nel numero reale più piccolo positivo, in modo tale da poterne calcolare il logaritmo.

La funzione get\_emission\_p restituisce la probabilità di emissione presente nel dizionario emission\_probabilities oppure la calcola in base alla strategia di smoothing utilizzata.

```
for s,tag in enumerate(possible_tags):
    transition_p = transition_matrix.loc['START',tag]
    emission_p = get_emission_p(emission_probabilities, sentence_tokens[0], tag,
count_word, smoothing_strategy, oneshot_words_tag_distribution, possible_tags)
    if transition_p == 0 : transition_p = np.finfo(float).tiny
    if emission_p == 0 : emission_p = np.finfo(float).tiny
    viterbi_matrix[s,0] = math.log(transition_p) + math.log(emission_p)
```

#### 2.Calcolo delle colonne successive

Vengono calcolati i valori di Viterbi degli stati delle colonne successive alla prima e man mano viene popolato il dizionario backpointer.

La funzione get\_max\_argmax\_value restituisce:

1. **max**\_: il prodotto massimo tra ogni valore di Viterbi della colonna precedente e la probabilità di transizione tra lo stato a cui fa riferimento e lo stato corrente.

$$\max_{s'=1}^{N} viterbi[s', t-1] * a_{s',s}$$

2. **la riga s' che massimizza questo prodotto**. Inoltre, l'indice della riga viene memorizzato nel dizionario backpointer.

Il valore di Viterbi dello stato corrente è calcolato sommando **max**\_ con il logaritmo della probabilità di emissione del token rappresentato dalla colonna corrente.

```
for t in range(1,len(sentence_tokens)):
    backpointer_column = dict()
    for s, tag in enumerate(possible_tags):
        max_ , backpointer_column[s] = get_max_argmax_value(possible_tags,
    viterbi_matrix, transition_matrix, t, s)
        emission_p = get_emission_p(emission_probabilities, sentence_tokens[t],
    tag, count_word, smoothing_strategy, oneshot_words_tag_distribution,
    possible_tags)
    if emission_p == 0: emission_p = np.finfo(float).tiny
        viterbi_matrix[s,t] = max_ + math.log(emission_p)
    backpointer[t] = backpointer_column
```

#### 3.Step finale

Calcolo del best\_path\_pointer, ovvero il <u>riferimento alla riga della cella dell'ultima colonna che massimizza la probabilità della sequenza di tag migliore</u>. Quest'ultima è calcolata moltiplicando il valore di Viterbi della cella in questione con la **probabilità di transizione finale**.

```
max_ = -sys.maxsize
best_path_pointer = None
for s,tag in enumerate(possible_tags):
end_transition = transition_matrix.loc[tag,'END']
if end_transition == 0: end_transition = np.finfo(float).tiny
val = viterbi_matrix[s,len(sentence_tokens) - 1] + math.log(end_transition)
if val >= max_: max_ = val ; best_path_pointer = s
```

$$\max_{s=1}^{N} viterbi[s,T] * a_{s,q_F}$$

## 4.Backtracking (riga 50 del codice)

Partendo dal **best\_path\_pointer** viene srotolata la sequenza dei tag più probabile sfruttando il dizionario backpointer. Attraverso la seguente immagine possiamo intuire facilmente il processo:

```
best_path_pointer = 11
```

```
.
23: {0: 6, 1: 6, 2: 6, 3: 6, 4: 6, 5: 6, 6: 6}
7: 6, 8: 6, 9: 6, 10: 6, 11: 6, 12: 6,
13: 6, 14: 6},
24: {0: 6, 1: 6, 2: 6, 3: 6, 4: 6, 5: 6, 6: 6, 7: 6, 8: 6, 9: 6, 10: 6, 11: 6, 12: 6,
13: 6, 14: 6},
25: {0: 6, 1: 6, 2: 6, 3: 6, 4: 6, 5: 6, 6: 6, 7: 6, 8: 6, 9: 6, 10: 6, 11: 6, 12: 6,
13: 6, 14: 6},
26: {0: 0, 1: 0, 2: 0, 3: 0, 4: 0, 5: 0, 6: 0, 7: 0, 8: 0, 9: 0, 10: 0, 11: 0) 12: 0,
13: 0, 14: 0}}
```

# 3. Strategie di Smoothing

Sono state utilizzate diverse strategie di **smoothing** per le parole sconosciute:

- P(unk|NOUN) = 1: se il tag preso in considerazione per la parola sconosciuta è **NOUN** allora la probabilità di emissione sarà uguale a 1, ovvero si decide con estrema certezza che la parola sconosciuta è un **sostantivo**
- P(unk|NOUN) = P(unk|VERB) = 0.5: se il tag preso in considerazione per la parola sconosciuta e **NOUN** o **VERB** allora la probabilità di emissione sarà uguale a 0.5, ovvero la parola può essere un sostantivo o un verbo con la stessa probabilità
- $P(unk|t_i) = \frac{1}{|POSTAGs|}$ : la proabilità di emissione per qualsiasi tag sarà 1 diviso la cardinalità dell'insieme dei tag possibili, ovvero una parola sconosciuta può essere un **qualsiasi tag** della lista dei tag possibili con la stessa probabilità
- **Statistica PoS sul development set**: la probabilità di emissione viene calcolata sulla base di una distribuzione estratta da un determinato corpus. In pratica, la distribuzione è relativa ai tag associati alle parole che compaiono **una sola volta** nel dev-set.

Il procedimento consiste nel collezionare tutte le parole che compaiono una sola volta (**one shot word**) e verificare quante di queste sono NOUN, VERB, ADJ, ecc...

Come tutte le distribuzioni, quella calcolata assocerà ad ogni tag una probabilità e la somma di tutte le probabilità sarà 1. Nell'algoritmo compute\_oneshot\_words\_distributions nel file learning2.py, che calcola e restituisce la distribuzione, verranno considerati anche i tag non associati a nessuna one shot word presente nel dev-set; a questi verrà assegnata probabilità nulla.

Ad esempio, la distribuzione di probabilità per i tag relativi alle one shot word del dev-set del Latino sarà definita come una lista di coppie (tag, probabilità):

```
[('NOUN', 0.19340159271899887), ('PROPN', 0.38680318543799774), ('VERB', 0.229806598407281), ('ADJ', 0.080773606370876), ('CCONJ', 0.0011376564277588168), ('DET', 0.051194539249146756), ('NUM', 0.012514220705346985), ('ADP', 0.009101251422070534), ('ADV', 0.022753128555176336), ('PRON', 0.004550625711035267), ('AUX', 0.005688282138794084), ('SCONJ', 0.0011376564277588168), ('PART', 0.0011376564277588168), ('PUNCT', 0), ('X', 0)]
```

E' possibile notare un'alta percentuale di nomi propri (PROPN) perchè rappresentano quella parte del discorso che viene menzionata più raramente, specialmente se trattiamo un corpus con sentence che non sono in relazione contestuale tra di loro.

# 4. Valutazione del Sistema e Analisi degli errori

# 4.1 Risultati per il Latino

#### **Baseline**

L'algoritmo di Baseline ottiene **ottimi** risultati. I vantaggi di questo algoritmo sono l'estrema semplicità e la velocità di esecuzione, infatti risulta essere quasi istantaneo. Gli errori più comuni sono principalmente relativi alla valutazione dei **nomi propri** e **verbi**.

```
Accuratezza: 95.39 %
Conteggi errori: {'VERB': 248, 'PROPN': 471, 'ADV': 56, 'DET': 142, 'NUM': 34, 'ADJ': 83, 'NOUN': 15, 'CCONJ': 35, 'ADP': 6, 'SCONJ': 15, 'AUX': 3, 'PRON': 2}
Tempo di esecuzione: 0.52 sec
```

#### Viterbi

L'algoritmo di Viterbi ottiene risultati leggermente migliori del Baseline per quanto riguarda le tre strategie di smoothing, ma in tutti i casi il tempo di esecuzione risulta essere estremamente più lungo. Anche in questo caso gli errori più comuni riguardano la valutazione dei **nomi propri**; questo accade perchè i nomi propri sono quelli che appaiono più raramente nel train set e molti di questi sono trattati come parole sconosciute; di conseguenza risulta essere determinante "l'azzardo" effettuato dalla strategia di smoothing utilizzata.

Notiamo che le strategie di smoothing sono fondamentali perchè per loro natura tendono a far diminuire gli errori riguardo i tag su cui azzardano; nel primo smoothing, il conteggio degli errori relativi ai nomi è molto basso, nel secondo invece aumenta ma diminuisce quello relativo ai verbi.

La strategia di smoothing migliore risulta essere l'ultima, perchè viene ottenuta un'approssimazione molto vicina a quella che sarebbe la distribuzione di probabilità dei tag relativi alle parole sconosciute.

```
Tipologia di smoothing: UNKNOWN_NAME
Accuratezza: 95.98 %
Conteggi errori: {'VERB': 200, 'PROPN': 471, 'ADV': 54, 'NUM': 15, 'PRON': 32, 'DET': 25, 'ADJ': 82, 'NOUN': 19, 'AUX': 31, 'CCONJ': 21, 'PUNCT': 5, 'SCONJ': 12, 'ADP': 2}
Tempo di esecuzione: 61.71 sec
```

```
Tipologia di smoothing: UNKNOWN_NAME_VERB
Accuratezza: 96.22 %
Conteggi errori: {'VERB': 121, 'PROPN': 471, 'ADV': 49, 'NUM': 15, 'PRON': 22, 'DET': 25, 'ADJ': 82, 'NOUN': 56, 'AUX': 29, 'CCONJ': 21, 'PUNCT': 4, 'SCONJ': 12, 'ADP': 2}
Tempo di esecuzione: 59.65 sec
```

```
Tipologia di smoothing: UNKNOWN_ALL
Accuratezza: 96.42 %
Conteggi errori: {'VERB': 148, 'PROPN': 387, 'ADV': 48, 'NUM': 15, 'PRON': 22, 'DET': 23, 'ADJ': 82, 'NOUN': 69, 'AUX': 29, 'CCONJ': 21, 'PUNCT': 5, 'SCONJ': 12, 'ADP': 2}
Tempo di esecuzione: 58.13 sec
```

```
Tipologia di smoothing: UNKNOWN_DISTRIBUTION_ONESHOT_WORDS
Accuratezza: 97.22 %
Conteggi errori: {'VERB': 159, 'PROPN': 184, 'ADV': 49, 'NUM': 15, 'PRON': 22, 'DET': 25, 'ADJ': 82, 'NOUN': 66, 'AUX': 29, 'CCONJ': 21, 'PUNCT': 4, 'SCONJ': 12, 'ADP': 2}
Tempo di esecuzione: 58.18 sec
```

## 4.2 Risultati per il Greco

#### **Baseline**

Le performance per il Greco calano di molto. L'algoritmo di Baseline ha un'accuratezza del 73.5 % e gli errori più comuni riguardano i **verbi**, gli **avverbi** e i **nomi**.

```
Accuratezza: 73.53 %

Conteggi errori: {'VERB': 1978, 'ADV': 1823, 'PRON': 462, 'ADJ': 965, 'CCONJ': 133, 'DET': 88, 'SCONJ': 25, 'NOUN': 50, 'ADP': 16, 'NUM': 1, 'PUNCT': 3, 'INTJ': 3, 'X': 1}

Tempo di esecuzione: 0.49 sec
```

#### Viterbi

Anche per il Greco, la strategia di smoothing migliore risulta essere quella relativa alla distribuzione di probabilità dei tag relativi alle parole che compaiono una sola volta nel dev set.

```
Tipologia di smoothing: UNKNOWN_NAME
Accuratezza: 73.60 %
Conteggi errori: {'VERB': 1974, 'ADV': 1710, 'PRON': 527, 'ADJ': 959, 'CCONJ': 125, 'DET': 103, 'SCONJ': 61, 'NOUN': 47, 'ADP': 18, 'PUNCT': 4, 'NUM': 1, 'INTJ': 3, 'X': 1}
Tempo di esecuzione: 43.34 sec
```

```
Tipologia di smoothing: UNKNOWN_NAME_VERB
Accuratezza: 76.43 %
Conteggi errori: {'VERB': 725, 'ADV': 1646, 'PRON': 501, 'NOUN': 780, 'ADJ': 966, 'CCONJ': 130, 'DET': 112, 'SCONJ': 52, 'ADP': 18, 'PUNCT': 4, 'NUM': 1, 'INTJ': 3, 'X': 1}
Tempo di esecuzione: 43.31 sec
```

```
Tipologia di smoothing: UNKNOWN_ALL
Accuratezza: 73.94 %

Conteggi errori: {'VERB': 974, 'ADV': 1638, 'PRON': 467, 'NOUN': 1079, 'ADJ': 970, 'CCONJ': 131, 'DET': 129, 'SCONJ': 50, 'ADP': 16, 'PUNCT': 3, 'NUM': 1, 'INTJ': 3, 'X': 1}

Tempo di esecuzione: 43.07 sec
```

```
Tipologia di smoothing: UNKNOWN_DISTRIBUTION_ONESHOT_WORDS
Accuratezza: 76.25 %
Conteggi errori: {'ADV': 1633, 'PRON': 494, 'VERB': 380, 'NOUN': 1172, 'ADJ': 968, 'CCONJ': 131, 'DET': 123, 'SCONJ': 49, 'ADP': 18, 'PUNCT': 4, 'NUM': 1, 'INTJ': 3, 'X': 1}
Tempo di esecuzione: 43.81 sec
```