

# Part-of-Speech (POS) Tagging

## 1. Perceptron Algorithm.

- Given word sequence  $X^{1 \dots S}$ , POS tagging  $Y^{1 \dots S}$ , feature  $\phi$   
Initialize  $w \leftarrow 0$

Repeat till converge:

For sentence  $s = 1$  to  $S$  do

$$V^s = \arg \max_Y \sum_{j=1}^I w_j \sum_{i=1}^{N_i} \phi_i(X^s, Y, j) \Rightarrow \text{Using Viterbi Decoding}$$

if  $V^s \neq Y^s$  do

$$f(X^s, Y^s) \leftarrow \sum_{i=1}^{N_i} \phi(X^s, Y^s, i)$$

$$f(X^s, V^s) \leftarrow \sum_{i=1}^{N_i} \phi(X^s, V^s, i)$$

$$\hat{w} += f(X^s, Y^s) - f(X^s, V^s)$$

- Viterbi Decoding.

~~$V^s$~~  Allocate  $\delta$  with size  $(|T|N)$        $T$ -#tags

For  $i = 1$  to  $N$  do

For  $t \in T$  do

$$\delta(i, t) = -\infty$$

For  $t' \in T$  do

$$\delta(i, t) = \max \left\{ \delta(i, t), \delta(i-1, t') + \sum_{j=1}^I w_j \phi_j(X, Y, i) \right\}$$

return  $\delta$

2. HMM:  $P(y_1^N, x_1^N) = \prod_{i=1}^N P(y_i | y_{i-1}) P(x_i | y_i)$   $y$ -hidden state,  $x$ -observation.

- Viterbi Decoding: find the best tagging (hidden) sequence, given sentences of words.

$$y_1^* = \underset{y_1^N}{\operatorname{argmax}} P(y_1^N | x_1^N) \Leftrightarrow P(y_1^N, x_1^N)$$

Allocate  $\delta$  with size  $|T|N$ ,  $\delta = -\infty$

For  $i = 1$  to  $N$  do

For  $t = 1$  to  $T$  do

For  $t' = 1$  to  $T$  do

$$\delta(i, t) = \max \{ \delta(i, t), \delta(i-1, t') P(t | t') P(x_i | t) \}$$

return  $\delta$

- Minimum Bayes Risk:  $\sum_{i=1}^N \max_{\hat{y}_i} P(\hat{y}_i = y_i | x_1^N)$ .

Forward-Backward Algorithm.

define:  $\alpha(i, t) = P(x_i^i, y_i = t)$

$$\beta(i, t) = P(x_{i+1}^N | y_i = t)$$

Allocate  $\alpha$  with size  $|T|N$ ,  $\alpha = 0$ ,  $\alpha(1, \text{START}) = 1$

For  $i = 1$  to  $N$  do

For  $t = 1$  to  $T$  do

For  $t' = 1$  to  $T$  do

~~$$\alpha(i, t) = \alpha(i, t) +$$~~

$$\alpha(i, t) += P(y_i = t | y_{i-1} = t') P(x_i | y_i) \alpha(i-1, t')$$

Allocate  $\beta$  with size  $|T|N$ ,  $\beta = 0$ ,  $\beta(N, :) = 1$ .

For  $i = N$  to  $1$  do

For  $t = 1$  to  $T$  do

For  $t' = 1$  to  $T$  do

$$\beta(i, t) += P(y_{i+1} = t' | y_i = t) P(x_{i+1} | y_{i+1}) \beta(i+1, t')$$

Best tagging sequence:

$$\left[ \underset{\hat{y}_i}{\operatorname{argmax}} \alpha(i, \hat{y}_i) \beta(i, \hat{y}_i) \right]_{i=1}^N$$

- EM Algorithm for parameters learning.

First compute  $\alpha$  and  $\beta$  for all sequences in data.

E-step: expected count:

$$\left\{ \begin{aligned} ec(t, x) &= \sum_{i, s.t. X_i = x} P(Z_i = t | X_i^N) = \frac{\sum_i \frac{1}{Z} \alpha(i, t) \beta(i, t)}{\sum_i \frac{1}{Z} \alpha(i, t) P(t|t) P(x_{i+1}|t') \beta(i+1, t')} \\ ec(t, t') &= \sum_i P(Z_i = t \wedge Z_{i+1} = t' | X_i^N) \\ &= \frac{\sum_i \frac{1}{Z} \alpha(i, t) P(t'|t) P(x_{i+1}|t') \beta(i+1, t')}{\sum_i \frac{1}{Z} \alpha(i, t) P(t'|t) P(x_{i+1}|t') \beta(i+1, t')} \end{aligned} \right.$$

M-step:

$$O_{t,x} = \frac{ec(t, x)}{\sum_{x'} ec(t, x')}$$

$$A_{t,t'} = \frac{ec(t, t')}{\sum_{t''} ec(t, t')}$$



3. Conditional Random Field (CRF):  $P(Y|X) = \frac{1}{Z} e^{\sum_j w_j f_j(Y, X)}$ . (max entropy).

- Learn the weights of CRF:

First calculate  $\alpha$  and  $\beta$  using sentences in data.

where the update rule of  $\alpha$  is modified as:

$$\alpha(i, t) += \alpha(i-1, t') e^{\sum_j w_j f_j(x, y_{i-1}, y_i)}$$

$$\beta(i, t) += \beta(i+1, t') e^{\sum_j w_j f_j(x, y_i, y_{i+1})}$$

Update  $w$  through gradient descent: (CRF has concave loss function)

$$\frac{\partial L}{\partial w} = f(Y, X) - \mathbb{E}_{P_w(Y|X)}[f]$$

where  $\mathbb{E}_{P_w(Y|X)}[f_j] = \sum_i \sum_{t, t'} f_j(Y_i=t, Y_{i+1}=t', X_i^N) \underbrace{\alpha(i, t) \beta(i+1, t') \frac{1}{Z} e^{w^T f(Y_i, Y_{i+1}, X_i^N)}}_{P(Y_i=t, Y_{i+1}=t' | X_i^N)}$

$$w += \eta \frac{\partial L}{\partial w}$$

- Latent CRF:  $P(Z, Y|X) = \frac{1}{Z(X)} e^{\sum_j w_j f_j(Z, Y, X)}$

$$\frac{\partial L}{\partial w} = \mathbb{E}_{P(Z, Y, X)}[f] - \mathbb{E}_{P(Z, Y|X)}[f]$$

4. Structured SVM: max margin as a metric for structured prediction.

$$\min \frac{1}{2} \|w\|^2 + C \sum_n \xi_n$$

$$\text{s.t. } w^T (f(x_n, y_n) - f(x_n, y)) \geq L(y_n, y) - \xi_n \quad \forall n, y \rightarrow \text{to all possible } y$$

where  $L(y_n, y) = \sum_{i=1}^{|x_n|} \mathbb{I}(y_{n,i} \neq y_i)$  is loss function.

- Training: (compared with perceptron for tagging)

$$\text{update rule: } w \leftarrow w - \eta \left( \frac{w}{\|w\|} \cdot C(f(x_n, y_n) - f(x_n, \hat{y})) \right)$$

$$\text{where } \hat{y} = \arg \max_y L(y_n, y) - w^T (f(x_n, y_n) - f(x_n, y))$$

- Solve  $\hat{y}$  through Loss Augmented Decoding:

For  $i = 1$  to  $|x_n|$  do

For  $t = 1$  to  $T$  do

For  $t' = 1$  to  $T$  do

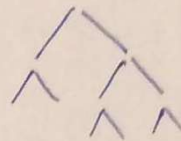
$$s(i, t) = \max \{ s(i, t), s(i-1, t') + l(x_n, y_{n,i}, t) - w^T (f(x_n, y_{n,i-1}, y_{n,i}) - f(x_n, t', t)) \}$$

# Parsing

- Key Points:
- Context Free Grammar (CFG) / Probabilistic CFG (PCFG)
  - Viterbi Decoding for parsing
  - Posterior Decoding for parsing (Inside-Outside Algorithm)

## 1. CFG.

- Parsing Problem: to define a syntactic structure (parsing tree)
- CFG:  $G(S, N, T, R)$ :  $S$  - start symbol,  $N$  - non-terminals,  $T$  - terminals,  $R$  - rules ( $X \rightarrow Y_1 Y_2 \dots Y_n$ , for  $n \geq 0$ ,  $X \in N$ ,  $Y_i \in (N \cup T)$ )
- PCFG:  $\sum_B P(A \rightarrow B) = 1 \cdot \forall A$   
 $P(\text{derivation}) = \prod_{r \in \text{tree}} P(r) \Rightarrow \text{probability of a tree: } p(t) = \prod_{i=1}^n q_i(\alpha_i \rightarrow \beta_i)$
- Hmm is a special case of PCFG:



## 2. Viterbi Decoding: With probabilities of all rules in CFG, find the best parsing tree

- $\text{tree}^* = \arg \max_{\text{tree}} P(\text{sentence, tree}) = \arg \max_{\text{tree}} \prod_{r \in \text{tree}} P(r)$
- Dynamic Programming:  $O(N^3 |R|^3)$   
define  $\delta(i, j, A)$  = maximum prob. of a constituent with non-terminal  $A$  spanning words  $i \dots j$  inclusive
- Chomsky Normal Form (CNF):  
 $\begin{cases} A \rightarrow BC \\ A \rightarrow a \end{cases}$   
every CFG can be transformed into CNF

Allocate  $\delta$  with size  $(|R| \cdot N^2)$ , all  $\delta = -\infty$

For  $i = 1$  to  $N$  do

For all  $X \in N$

$$\delta(i, i, X) = \begin{cases} P(X \rightarrow x_i) & \text{if } X \rightarrow x_i \in R \\ 0 & \text{otherwise} \end{cases}$$

For span = 2 to  $N$  do

For  $i = 1$  to  $N+1 - \text{span}$  do

$j = i + \text{span}$

For  $k = i$  to  $j$  do

For  $A, B, C \in N$  do

$$\delta(i, j, A) = \max \{ \delta(i, j, A), P(A \rightarrow BC) \delta(i, k, B) \delta(k, j, C) \}$$

return  $\delta(1, N+1, S)$



3. Posterior Decoding: despite of maximizing tree prob., also maximize the correctness of each node (Minimum Bayesian Risk, MBR).

$$\max_t \sum_{\text{net}} P(\text{net}) \Rightarrow P(\text{net}) = \frac{\alpha(i, j, A) \beta(i, j, A)}{\beta(1, N+1, S)}$$

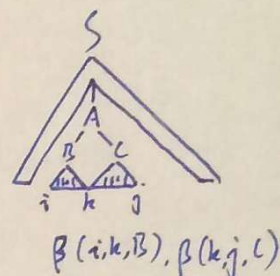
• Inside-Outside Algorithm.

define:  $\alpha(i, k, X) = P(S \rightarrow^* x_1 \dots x_i X x_{i+1} \dots x_{N+1})$

$$\alpha(i, j, X) = P(S \rightarrow^* x_1 \dots x_i X x_{i+1} \dots x_{N+1})$$

$$\beta(i, k, X) = P(X \rightarrow^* x_i \dots x_k)$$

$$S(A \rightarrow BC, i, j, k) = \frac{1}{Z} \alpha(i, j, A) \beta(i, k, B) \beta(k, j, C) P(A \rightarrow BC)$$



Allocate  $\alpha$  and  $\beta$  with size  $|R| N^2$

$$\alpha(1, N, S) = 1$$

# calculate  $\beta$  first.

For ~~span~~  $i = 1$  to  $N$  do

For all  $X \in N$

$$\beta(i, i+1, X) = \begin{cases} P(X \rightarrow x_i), & \text{if } X \rightarrow x_i \text{ in } R \\ 0 & \text{otherwise.} \end{cases}$$

For span = 2 to  $N$  do

For  $i = 1$  to  $N+1$ -span do

$$j = i + \text{span}$$

For  $k = i$  to  $j$  do

For  $A, B, C \in N$  do

$$\beta(i, k, A) += \beta(i, k, B) \beta(k, j, C) P(A \rightarrow BC)$$

# then calculate  $\alpha$

For span =  $N$  to 2 do

For  $i = 1$  to  $N+1$ -span do

$$j = i + \text{span}$$

For  $k = i$  to  $j$  do

For  $A, B, C \in N$  do

$$\alpha(i, k, B) += \alpha(i, j, A) \beta(k, j, C) P(A \rightarrow BC)$$

$$\alpha(k, j, C) += \alpha(i, j, A) \beta(i, k, B) P(A \rightarrow BC)$$

return  $\alpha$  and  $\beta$

- Calculate  $S$  with  $\alpha$  and  $\beta$ : build the most likely tree, but need to have a valid tree structure

Allocate  $S$  with size  $|R|N^2$ ,  $S = -\infty$

For span = 2 to  $N$  do

For  $i = 1$  to  $N + 1 - \text{span}$  do

$j = i + \text{span}$

For  $k = i$  to  $j$  do

For  $A, B, C \in N$  do

$$S(i, j, A) = \max \{ S(i, j, A), S(i, k, B) + S(k, j, C) + \alpha(i, j, A) \beta(i, j, A) \}$$

return  $S$



# Machine Translation

- Key Points:
- Word Alignment (IBM Models)
  - Phrasal Based Translation
  - Synchronous CFG

## 1. Word Alignment

Task: French  $\rightarrow$  English

- Noisy Channel Model

$$p(e|f) = \frac{p(f|e)p(e)}{p(f)}$$

$p(e)$ : language model

$p(f|e)$ : translation model (backwards!)

- IBM Model 1

alignment sequence:  $a: \{a_1, a_2, \dots, a_m\}$   $m$  - french sentence length.  $l$  - Eng. sentence length

$$p(f, a|e, m) = p(a|e, m) p(f|a, e, m) = p(a|e, m) p(f|e_{a_j})$$

alignment prob.                  translation prob.

- Assumption 1:  $p(a|e, m) = \frac{1}{(l+1)^m}$  (every alignment is equally probable)  
constant.

- Assumption 2: Each ~~target~~ <sup>source</sup> word corresponds to only one target word.  
 $p(f_j|e_{a_j})$  — learned using EM.

$$\text{IBM Model 1: } p(f, a|e, m) = \frac{c}{(l+1)^m} \prod_j p(f_j|e_{a_j})$$

- IBM Model 2

- Introduce distortion parameters:  $q_e(i|j, l, m)$   
index of Eng.                  index of French.

$$p(a|e, m) = \prod_{j=1}^m q_e(a_j|j, l, m)$$

$$\text{IBM Model 2: } p(f, a|e, m) = \prod_{j=1}^m p(f_j|e_{a_j}) q_e(a_j|j, l, m)$$

- HMM Model

$$p(f, a|e, m) = \prod_{j=1}^m p(f_j|e_{a_j}) p(a_j|a_{j-1}, l, m)$$



• EM for IBM Model 2.

Initial  $q_e()$ ,  $p_i()$  to be random values

Allocate  $\delta$  with size  $|E||F|$

E-step: For all sentence:

For  $i = 1$  to  $m$  do.

For  $j = 0$  to  $l$  do // 0 for NULL symbol.

$$\text{count}(e_j, f_i) += p(f_i | e_j) q_e(j | i, l, m) \frac{1}{Z}$$

$$\text{count}(e_j) += \dots$$

$$\text{count}(j, i, l, m) += \dots$$

$$\text{count}(i, l, m) += \dots$$

~~Normalise~~

$$M\text{-step: } p(f_i | e_j) = \frac{\text{count}(e_j, f_i)}{\text{count}(e_j)} \quad q_e(j | i, l, m) = \frac{\text{count}(j, i, l, m)}{\text{count}(i, l, m)}$$

## 2. Phrased Based Translation

### ① Alignment Matrix

Step 1: train IBM Model 2 for  $p(f|e)$ , find best alignment:  $a_1^* = \arg\max_a p(f, a|e, m)$

Step 2: train - - - for  $p(e|f)$ , find best alignment:  $a_2^* = \arg\max_a p(e, a|f, l)$

Step 3: Intersection of two most likely alignment.

Step 4: Growth algorithm to get alignment matrix.

### ② Extract Phrase Pairs from Alignment Matrix.

• Pair  $(e, f)$  is consistent if:

- (1) at least one word in  $e$  aligned to a word in  $f$
- (2) no words in  $f$  aligned to words outside  $e$ .
- (3) no words in  $e$  aligned to words outside  $f$ .

	Mary did not
Maria	?
has	o o

• Probabilities of Phrase Pairs

$$p(f|e) = \frac{\text{count}(f, e)}{\text{count}(e)}$$

### ③ Decoding with Phrased Based Model:

• Do not consider reordering: only phrase model (transition prob.) and language model.

Dynamic Programming: state =  $[i, e, \underline{e'}]$  last two english words, for tri-gram language model

For  $i = 1$  to  $m$  do

For  $j = 1$  to  $i$  do // determine the length of last phrase

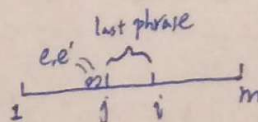
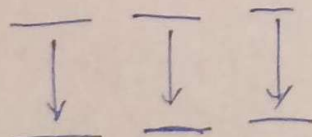
For phrase  $ph$  s.t. lexicon  $(p, f_j^i)$  do

For  $e, e'$  do

$$\delta(i, e, e') \max = \delta(j, e, e') + \log p(ph|f_j^i) + P_{lm}(ph|e, e')$$

French

English





- Consider reordering and distortion parameters.

Beam Search: state:  $(e, e', b, r, a)$  score initial state  $(*, *, 0^m, 0, 0)$   
 bit-string end-point of last phrase  
 indicating covered position

functions:  $\phi(q)$ : possible states  $\left\{ \begin{array}{l} \text{must not overlap with } b \\ \text{distortion limit must not be violated.} \end{array} \right.$

- $\text{next}(q, p)$ : combining state  $q$  with phrase  $p$

- $eq(q_1, q_2)$ : return true if  $e=e_1, e'=e'_2, b=b_2, r=r_2$  (ignoring scores)

• Add  $(Q, \hat{q}_e, q_e, p)$ : If  $eq(\hat{q}_e, \hat{q})$ :  
If  $\alpha(\hat{q}_e) > \alpha(\hat{q})$ :  
 $\mathcal{Q} = \{\hat{q}_e\} \cup \mathcal{Q} \setminus \{\hat{q}_e\}$

Else :

$$\mathcal{Q} = \mathcal{Q} \cup \{\hat{q}_k\}$$

- $\text{beam}(\mathcal{Q}) = \{q_i \in \mathcal{Q} : a(q_i) \geq a^* - \beta\}$

$\beta z_0$ : beam-width parameter

$$a^* = \max_{q \in Q} a(q)$$

Initialize  $Q_0 = \{q_0\}$ ,  $Q_i = \{\phi\}$

For  $i = 0$  to  $m-1$

For each state  $q_i \in \text{beam}(Q_i)$  do

For each phrase  $p \in \text{ph}(q_e)$  do

$$\hat{q}_k = \text{next}(q_k, p)$$

Add ( $Q_i, \hat{q}_i, q_i, p$ ) (where  $i = \text{len}(\hat{q}_i)$ )

return the highest score state in  $Q_n$ .

score: 
$$a(\hat{q}_e) = a(q_e) + \underbrace{\log p(\text{fle})}_{\text{(phrase model)}} + \underbrace{\log p_{2m}(\hat{q}_e)}_{\text{(language model)}} + \eta \times |j+1 - r(q_e)|.$$
 (distortion).

French

English

