

10/28

**Review**

\* ML estimation for incomplete data


Examples  $t = 1, 2, \dots, T$ Visible nodes  $V^{(t)}$ 

\* EM algorithm

E-step: compute posteriors  $P(X_i = x, \pi_i = \pi | V^{(t)})$  inference

M-step: update CPTs

$$P(X_i = x | \pi_i = \pi) \leftarrow \frac{\sum_{\pi} P(X_i = x, \pi_i = \pi | V^{(t)})}{\sum_{\pi} P(\pi_i = \pi | V^{(t)})}$$

Converges to local maximum of  $\mathcal{L} = \sum_t \log P(V^{(t)})$ **Example**

 A and C are observed. B is hidden.

\* Posterior probability

$$P(B=b | A=a, C=c) = \frac{P(C=c | B=b, A=a) P(B=b | A=a)}{\sum_{b'} P(C=c | B=b', A=a) P(B=b' | A=a)} \quad \text{Bayes rule}$$

$$P(b | a, c) = \frac{P(c | b) P(b | a)}{\sum_{b'} P(c | b') P(b' | a)}$$

\* Incomplete data set  $\{ (a_t, c_t) \}_{t=1}^T$  (I.I.D.)

$$\begin{aligned} \text{Log-likelihood } \mathcal{L} &= \sum_t \log P(a_t, c_t) \\ &= \sum_t \log \sum_b P(a_t, b, c_t) \\ &= \sum_t \log \sum_b [P(a_t) P(b | a_t) P(c_t | b)] \end{aligned}$$

M-step for this example:

$$P(B=b | A=a) \leftarrow \frac{\sum_t P(A=a, B=b | A=a_t, C=c_t)}{\sum_t P(A=a | A=a_t, C=c_t)}$$

$$\text{Simplify RHS} = \frac{\sum_t I(a, a_t) P(b | a_t, c_t)}{\sum_t I(a, a_t)}$$

$$P(C=c | B=b) \leftarrow \frac{\sum_t P(c, b | A=a_t, C=c_t)}{\sum_t P(b | A=a_t, C=c_t)}$$

$$\text{Simplify RHS} = \frac{\sum_t I(c, c_t) P(b | a_t, c_t)}{\sum_t P(b | a_t, c_t)}$$

Ex: Markov models of language.

\* Let  $w_l$  denote  $l^{\text{th}}$  word in text.How to model  $P(w_1, w_2, \dots, w_L)$ ?

Model	$P(\vec{w})$	ML estimate	DAG
unigram	$\prod_i P_i(w_{e_i})$	$P_i(w) = \frac{\text{count}(w)}{\sum_{w'} \text{count}(w')}$	$(w_1) \quad (w_2) \quad \dots \quad (w_L)$
bigram	$\prod_i P_2(w_{e_i}   w_{e_{i-1}})$	$P_2(w'   w) = \frac{\text{count}(w, w')}{\text{count}(w)}$	$(w_1) \rightarrow (w_2) \rightarrow \dots \rightarrow$

### \* Evaluating n-gram models

train on corpus A :  $P_1(\vec{w}) \leq P_2(\vec{w}) \leq P_3(\vec{w}) \leq \dots$

test on corpus B :  $P_2(\vec{w}) = 0$  if there are unseen bigrams.

$P_3(\vec{w}) = 0$  " " " " " trigrams.

### Word clustering

#### \* Alternative to bigram model



Replace by:



words  $w, w'$  observed  
cluster  $z$  hidden.

CPTs in BN:

$P(z | w)$  = prob that word  $w$  is mapped into cluster  $z$ .

$P(w' | z)$  = prob that word in cluster  $z$  is followed by word  $w'$ .

In cluster model:  $P(w' | w) = \sum_{z=1}^C P(w' | z) P(z | w)$ .

#### \* compact representation

# words =  $V$  (vocabulary size)

# clusters =  $C$  (clusters)

# parameters =  $2CV$  } reduce to unigram if  $C=1$

# bigrams =  $V^2$  } recover bigrams if  $C=V$ .

\* learning : how to estimate  $P(z | w)$  and  $P(w' | z)$ ?

E-step : 
$$P(z | w, w') = \frac{P(w' | z) P(z | w)}{\sum_{z=1}^C P(w' | z) P(z | w)}$$
 Bayes rule

M-step : update CPTs

$$P(z | w) \leftarrow \frac{\sum_i I(w, w_{e_i}) P(z | w_{e_i}, w_{e_{i+1}})}{\sum_i I(w, w_{e_i})}$$

$$P(w' | z) \leftarrow \frac{\sum_i I(w_{e_{i+1}}, w') P(z | w_{e_i}, w_{e_{i+1}})}{\sum_i P(z | w_{e_i}, w_{e_{i+1}})}$$

#### \* Experimental results

$V = 60,000$  word vocabulary.

$L = 80$  million word corpus of WSJ articles.

count( $w, w'$ ) is sparse : 99.8% elements are zero.

$$\text{count}(w, w') = \sum_i I(w, w_i) I(w', w_{i+1})$$

$C = 32$  model trained by EM.

$P(z|w)$  and  $P(w'|z)$  - approx 4 million parameters.

Converges in  $\sim 30$  iterations.

What clusters are discovered? For each word  $w$ , what is  $\max_z P(z|w)$ ?

### Linear Interpolation of Markov Models

$$P_M(w_l | w_{l-1}, w_{l-2}) = \lambda_1 \underset{\text{mixture model}}{P_1(w_l)} + \lambda_2 \underset{\text{unigram}}{P_2(w_l | w_{l-1})} + \lambda_3 \underset{\text{bigram}}{P_3(w_l | w_{l-1}, w_{l-2})}$$

$n$ -gram models are trained on corpus A.

How to estimate  $\lambda_i$  where  $\lambda_i \geq 0$  and  $\sum_{i=1}^3 \lambda_i = 1$ ?

#### \* Methodology

Train  $P_1, P_2, P_3$  on corpus A.

Fix  $P_1, P_2$ , and  $P_3$ .

Train  $\lambda_1, \lambda_2, \lambda_3$  on corpus C.

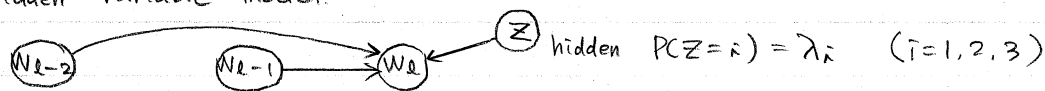
Estimate  $\lambda$  to maximize log-likelihood of corpus C.

- Don't estimate  $\lambda_i$ 's on corpus A :  $\lambda_1 = \lambda_2 = 0$   
 $\lambda_3 = 1$  } always favor trigram model.

- Test  $P_M = \sum_{i=1}^3 \lambda_i P_i$  on corpus B.

Don't estimate  $\lambda_i$ 's on corpus B : cheating.

#### \* Hidden variable model.



#### CPT for $w_l$

$$P(w_l | z, w_{l-1}, w_{l-2}) = \begin{cases} P_1(w_l) & \text{if } z=1 \\ P_2(w_l | w_{l-1}) & \text{if } z=2 \\ P_3(w_l | w_{l-1}, w_{l-2}) & \text{if } z=3 \end{cases}$$

$$\begin{aligned} \text{In this model: } P(w_l | w_{l-1}, w_{l-2}) &= \sum_{z=1}^3 P(w_l, z | w_{l-1}, w_{l-2}) \quad \text{marginalization} \\ &= \sum_{z=1}^3 P(z | w_{l-1}, w_{l-2}) P(w_l | z, w_{l-1}, w_{l-2}) \quad \text{product rule} \\ &= \sum_{z=1}^3 P(z) P(w_l | z, w_{l-1}, w_{l-2}) \quad \text{conditional independence} \\ &= \lambda_1 P_1(w_l) + \lambda_2 P_2(w_l | w_{l-1}) + \lambda_3 P_3(w_l | w_{l-1}, w_{l-2}) \end{aligned}$$

\* E-step

Compute posterior probability

$$P(Z=i | w_L, w_{L-1}, w_{L-2}) = \frac{P(w_L | Z=i, w_{L-1}, w_{L-2}) P(Z=i)}{\sum_j P(w_L | Z=j, w_{L-1}, w_{L-2}) P(Z=j)}$$

$$= \frac{\lambda_i P_i(w_L | Z=i, w_{L-1}, w_{L-2})}{\lambda_1 P_1(w_L) + \lambda_2 P_2(w_L | w_{L-1}) + \lambda_3 P_3(w_L | w_{L-1}, w_{L-2})}$$

\* M-step

Update parameters  $\lambda_i = P(Z=i)$

General EM:  $P(X_i = x | pa_i = \pi) \leftarrow \frac{\sum_{\mathcal{D}} P(X_i = x, pa_i = \pi | V^{(t)})}{\sum_{\mathcal{D}} \sum_{\pi} P(X_i = x, pa_i = \pi | V^{(t)})}$

Translation:  $i^{th}$  example  $\longleftrightarrow l^{th}$  word triplet

$X_i \longleftrightarrow \text{node } z$

$pa_i \longleftrightarrow \varnothing$  because  $z$  has no parents.

For mixture model:

$$P(Z=i) \leftarrow \frac{\sum_j P(Z=i | w_L, w_{L-1}, w_{L-2})}{\sum_{j=1}^3 \sum_k P(Z=j | w_L, w_{L-1}, w_{L-2})}$$

$$\lambda_i \leftarrow \frac{\sum_j P(Z=i | w_L, w_{L-1}, w_{L-2})}{L}$$

$L \leftarrow \# \text{ words in corpus.}$

\* Iterate EM: Guaranteed improvement of log-likelihood:

$$\mathcal{L}(\lambda_1, \lambda_2, \lambda_3) = \sum_{l=1}^L \log P_M(w_L | w_{L-1}, w_{L-2})$$

$$\lambda_1 P_1 + \lambda_2 P_2 + \lambda_3 P_3$$