

Review

* Learn CPTs from incomplete data.

$$P(X_i = x | p_{a_i} = \pi) \leftarrow \frac{\sum_{\pi} P(X_i = x, p_{a_i} = \pi | V^{(c)})}{\sum_{\pi} P(p_{a_i} = \pi | V^{(c)})}$$

EM algorithm

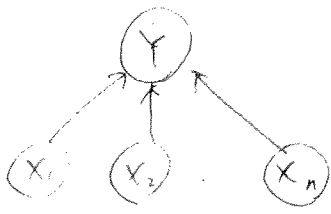
* Ex: Simple Chain with data $\{(a_t, c_t)\}_{t=1}^T$



$$P(b|a) \leftarrow \frac{\sum_{\pi} I(a, a_t) P(b|a_t, c_t)}{\sum_{\pi} I(a, a_t)}$$

Computed from Bayes rule

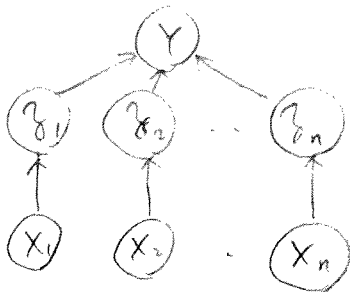
* Ex: noisy-OR with data $\{(\vec{x}_t, \vec{y}_t)\}_{t=1}^T$



$$P(Y=1 | \vec{x}) = 1 - \prod_{i=1}^n (1 - p_i)^{x_i}$$

* How to estimate noisy-OR parameters p_i ?

* Alternative network:



$$P(Z_i = 0 | X_i) = (1 - p_i)^{x_i}$$

$$P(Y=1 | \vec{x}) = 1 - \prod_{i=1}^n (1 - p_i)^{x_i} \quad \text{equivalent to noisy-OR}$$

* Conditional log-likelihood of data $\{(\vec{x}_t, y_t)\}_{t=1}^T$

$$\begin{aligned} \mathcal{L} &= \sum_t \log P(y_t | \vec{x}_t) \\ &= \sum_t [(1-y_t) \log P(y=0 | \vec{x}_t) + y_t \log P(y=1 | \vec{x}_t)] \\ &= \sum_t \left[(1-y_t) \log \prod_{i=1}^n (1-p_i)^{x_{it}} + y_t \log \left(1 - \prod_{i=1}^n (1-p_i)^{x_{it}} \right) \right] \\ &= \sum_t \left[(1-y_t) \sum_{i=1}^n x_{it} \log(1-p_i) + y_t \log \left(1 - \prod_{i=1}^n (1-p_i)^{x_{it}} \right) \right] \end{aligned}$$

↑ data
↑ data
↑ parameters

Note. Complicated, non-linear expression with respect to p_i !
EM to the rescue!

Shorthand: let T = total # examples.

let $T_i = \sum_{t=1}^T x_{it}$ (Count # times that $x_i=1$)

EM update rule:

$$P(z_i=1 | x_i=1) \leftarrow \frac{\sum_t P(z_i=1, x_i=1 | \vec{x}=\vec{x}_t, Y=y_t)}{\sum_t P(x_i=1 | \vec{x}=\vec{x}_t, Y=y_t)}$$

Simplify:

$$P(z_i=1 | x_i=1) \leftarrow \frac{\sum_t \mathbb{I}(x_{it}, 1) P(z_i=1 | \vec{x}_t, y_t)}{\sum_t \mathbb{I}(x_{it}, 1)}$$

Recall from last lecture:

$$P(z_i=1 | \vec{x}, y) = \frac{y p_i x_i}{1 - \prod_{i=1}^n (1-p_i)^{x_i}}$$

EM update rule:

$$p_i \leftarrow \frac{1}{T_i} \sum_t \mathbb{I}(x_{it}, 1) \left(\frac{y_t p_i x_{it}}{1 - \prod_{i=1}^n (1-p_i)^{x_{it}}} \right)$$

This update rule, applied in parallel to all $\{p_i\}_{i=1}^n$ will monotonically increase $\mathcal{L} = \sum_t \log P(y_t | \vec{x}_t)$

Markov Models of language

- * Let $w_l = l^{\text{th}}$ word in sentence.
How to model $P(w_1, w_2, \dots, w_L)$?
Shorthand $\vec{w} = (w_1, w_2, \dots, w_L)$.

Model	$P(\vec{w})$	ML estimate	DAG
unigram	$\prod_l P_1(w_l)$	$P_1(w) = \frac{\text{Count}(w)}{\sum_{w'} \text{Count}(w')}$	$(w_1) \quad (w_2) \quad \dots \quad (w_L)$
bigram	$\prod_l P_2(w_l w_{l-1})$	$P_2(w' w) = \frac{\text{Count}(w \rightarrow w')}{\text{Count}(w)}$	$(w_1) \rightarrow (w_2) \rightarrow \dots \rightarrow (w_L)$

* Evaluating n-grams

Train on corpus A: $P_1(\vec{w}) \leq P_2(\vec{w})$ on corpus A

Test on corpus B: $P_1(\vec{w}) \geq P_2(\vec{w})$ if $P_2(\vec{w}) = 0$ (these are unseen bigrams)

Linear interpretation

- * Also known as mixture model

$$P_m(w_l | w_{l-1}) = \lambda P_1(w_l) + (1-\lambda) P_2(w_l | w_{l-1})$$

How to estimate λ ?

* Methodology

Train P_1, P_2 on Corpus A

A = "training set"

Fix P_1 and P_2

B = "testing set"

estimate λ on corpus C

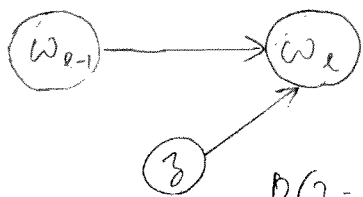
C = "development set"

Choose λ to maximize likelihood $\prod_l P_m(w_l | w_{l-1})$ on Corpus C.

- * Don't estimate λ on Corpus A: this would yield $\lambda = 1$ (always favor P_1)!

Do not estimate λ on Corpus B: cheating!

* Hidden Variable Model.



$$P(w_e | w_{e-1}, z) = \begin{cases} p_1(w_e) & \text{if } z=1 \\ p_2(w_e | w_{e-1}) & \text{if } z=2 \end{cases}$$

$$z \in \{1, 2\} \quad \left. \begin{array}{l} P(z=1) = \lambda \\ P(z=2) = 1-\lambda \end{array} \right\} \text{How to estimate on corpus } C?$$

Hidden Variable: z

Observed Variables: $\{(w_{e-1}, w_e)\}_{e=1}^L$ Corpus C

In This model:

$$\begin{aligned} P(w_e | w_{e-1}) &= \sum_{z=1}^2 P(w_e, z | w_{e-1}) \\ &= \sum_z P(w_e | z, w_{e-1}) P(z | w_{e-1}) \quad \text{product rule} \\ &= \sum_z P(w_e | z, w_{e-1}) P(z) \quad \text{conditional independence} \\ &= \lambda p_1(w_e) + (1-\lambda) p_2(w_e | w_{e-1}) \\ &= p_m(w_e | w_{e-1}) \quad \text{matches linear interpolation.} \end{aligned}$$

* E-step: Compute posterior probability \leftarrow conditional independence

$$P(z | w_{e-1}, w_e) = \frac{P(w_e | z, w_{e-1}) P(z | w_{e-1})}{P(w_e | w_{e-1})} \quad \text{Bayes rule}$$

$$P(z=1 | w_{e-1}, w_e) = \frac{\lambda p_1(w_e)}{\lambda p_1(w_e) + (1-\lambda) p_2(w_e | w_{e-1})}$$

* M-step:

General update rule EM:

$$P(x_i = x | p_{\alpha_i} = \pi) \leftarrow \frac{\sum_{\pi} P(x_i = x, p_{\alpha_i} = \pi | V^{(t)})}{\sum_{x'} \sum_{\pi} P(x_i = x', p_{\alpha_i} = \pi | V^{(t)})}$$

In this problem:

$$P(z=1) \leftarrow \frac{\sum_{\ell} P(z=1 | w_{\ell-1}, w_{\ell})}{\sum_{z'} \sum_{\ell} P(z'=z' | w_{\ell-1}, w_{\ell})}$$

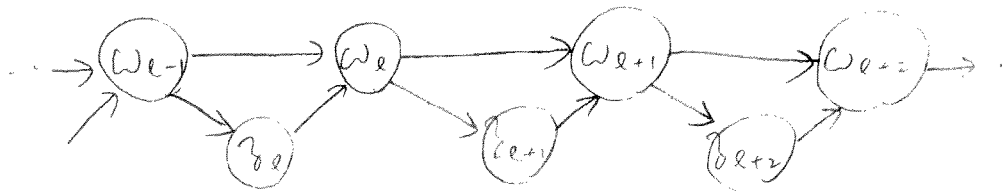
$$\lambda \leftarrow \frac{\sum_{\ell} P(z=1 | w_{\ell-1}, w_{\ell})}{L}$$

* Iterate EM:

Guaranteed improvement on corpus C of log-likelihood

$$L(\lambda) = \sum_{\ell} \log P_m(w_{\ell} | w_{\ell-1})$$

* In real-world application, smoothing parameter would depend on previous word $w_{\ell-1}$.



$$P(z=1 | w_{\ell-1}) = \lambda(w_{\ell-1}) \quad \lambda(w) \text{ depends on previous word}$$

$$P(z=2 | w_{\ell-1}) = 1 - \lambda(w_{\ell-1})$$

$$P(w_{\ell} | w_{\ell-1}) = \lambda(w_{\ell-1}) P_1(w_{\ell}) + (1 - \lambda(w_{\ell-1})) P_2(w_{\ell} | w_{\ell-1})$$

* EM used to estimate, say, $|V| \sim 10^5$ parameters.
 \uparrow
 Vocabulary size.