

Review

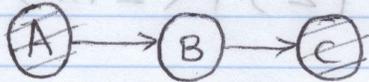
* Learning CPTs from incomplete data

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

E.M. Algorithm

* Example:

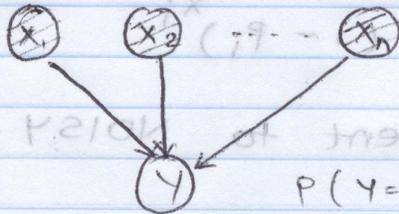
chain BN with data $\{(a_t, c_t)\}_{t=1}^T$



$$P(b|a) = \frac{\sum_{t=1}^T I(a, a_t) P(b|a_t, c_t)}{\sum_{t=1}^T I(a, a_t)}$$

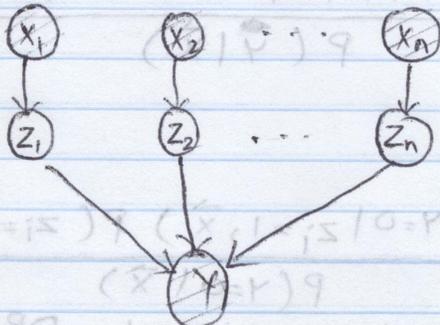
computed from Bayes rule and current CPTs

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



$$P(Y=1 | X_1, X_2, \dots, X_n) = 1 - \prod_{i=1}^n (1 - p_i)^{X_i}$$

Consider alternate BN:



$$P(Z_i = 1 | X_i = 1) = p_i$$

$$P(Z_i = 1 | X_i = 0) = 0$$

$$Y = \text{OR}(X_1, X_2, \dots, X_n)$$

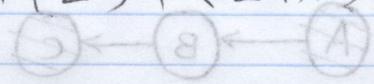
Equivalently: $P(z_i=0|x_i) = (1-p_i)^{x_i} = \begin{cases} 1-p_i & \text{if } x_i=1 \\ 1 & \text{if } x_i=0 \end{cases}$

What is $P(y=1|\vec{x})$ in this new model?

$$P(y=1|\vec{x}) = \sum_{\vec{z} \in \{0,1\}^n} P(y=1, \vec{z} | \vec{x}) \text{ marginalization}$$

$$= \sum_{\vec{z}} P(y=1 | \vec{z}, \vec{x}) P(\vec{z} | \vec{x}) \text{ product rule}$$

$$= \sum_{\vec{z}} P(y=1 | \vec{z}) P(\vec{z} | \vec{x}) \text{ conditional independence}$$



$$= \sum_{\vec{z} \neq \vec{0}} (1) P(\vec{z} | \vec{x}) + 0 \cdot P(\vec{z} = \vec{0} | \vec{x})$$

$$= 1 - P(\vec{z} = \vec{0} | \vec{x}) \text{ due to}$$

normalization

$$= 1 - \prod_{i=1}^n P(z_i=0|x_i) \text{ conditional independence}$$

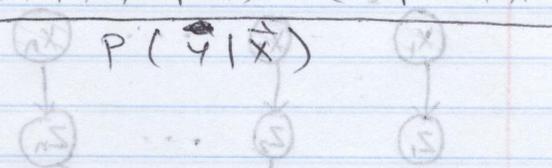
$$= 1 - \prod_{i=1}^n (1-p_i)^{x_i}$$

Equivalent to NOISY-OR!

1. Inference for EM algorithm:

compute posterior prob:

$$P(z_i=1 | \vec{x}, y) = \frac{P(y | \vec{x}, z_i=1) P(z_i=1 | \vec{x})}{P(y | \vec{x})}$$



$$P(z_i=1 | \vec{x}, y=0) = \frac{P(y=0 | z_i=1, \vec{x}) P(z_i=1 | \vec{x})}{P(y=0 | \vec{x})} = 0 \text{ due to logical-OR}$$

p_i if $x_i = 1$
 $1 - p_i$ if $x_i = 0$

$$P(z_i = 1 | \vec{x}, y = 1) = \frac{P(y = 1 | \vec{x}, z_i = 1) P(z_i = 1 | \vec{x})}{P(y = 1 | \vec{x})}$$

$$= \frac{(1 - p_i)^{x_i} p_i^{1 - x_i}}{\prod_{i=1}^n (1 - p_i)^{x_i} p_i^{1 - x_i}}$$

most combining:

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

* How good is model?

Compute conditional log-likelihood of data $\{(x_t, y_t)\}_{t=1}^T$

$$\mathcal{L} = \sum_{t=1}^T \log P(y_t | \vec{x}_t) = \sum_{t=1}^T \left[(1 - y_t) \log P(y = 0 | \vec{x}_t) + y_t \log P(y = 1 | \vec{x}_t) \right]$$

$$= \sum_{t=1}^T \left\{ (1 - y_t) \log \prod_{i=1}^n (1 - p_i)^{x_{it}} + y_t \log \left[1 - \prod_{i=1}^n (1 - p_i)^{x_{it}} \right] \right\}$$

data parameters to tune

Note: complicated, non-linear expression with respect to p_i !

Use EM to derive simple updates

Short hand: let $T =$ total # examples

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

$1 = x_{it} = 1$
 $0 = x_{it} = 0$

EM update rule: $P_i = (1 - P_i, \hat{x}_{it} | 1 = i, \hat{x}_{it})$

$$P(z_i = 1 | x_i = 1) \leftarrow \frac{\sum_{t=1}^T P(z_i = 1, x_i = 1 | \vec{x} = \vec{x}_t, y = y_t)}{\sum_{t=1}^T P(x_i = 1 | \vec{x} = \vec{x}_t, y = y_t)}$$

Simplify:

$$P(z_i = 1 | x_i = 1) \leftarrow \frac{\sum_{t=1}^T I(x_{it}, 1) P(z_i = 1 | \vec{x}_t, y_t)}{\sum_{t=1}^T I(x_{it}, 1)}$$

(computed from Bayes rule)

Final update:

$$P_i \leftarrow \left(\frac{1}{T_i} \right) \sum_{t=1}^T x_{it} \left(\frac{y_t P_i x_{it}}{1 - \prod_{j=1}^n (1 - P_j)^{x_{jt}}} \right)$$

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

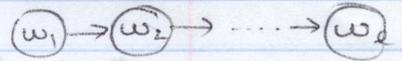
Markov models of language

* Let w_t be the t^{th} word in a sentence. How to model $P(w_1, w_2, \dots, w_L)$?
 Shorthand, $\vec{w} = (w_1, w_2, \dots, w_L)$.

Model $P(\vec{w})$? λ MLE estimate with DAG

unigram $\prod_e P_1(w_e)$ $P_1(w)$ $n \text{ count}(w)$ $(w_1) (w_2) \dots (w_n)$

bigram $\prod_e P_2(w_e | w_{e-1})$ $P_2(w' | w)$ $n \text{ count}(w \rightarrow w')$



Belief network

* Evaluating n -gram models

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 $P_2(\vec{w}) = 0$

no estimates of λ $\lambda = (1-s)$? $\{s, 1\}$ there are unseen bigrams
 $\lambda - 1 = (s = s)$?

Linear Interpolation

* Also known as mixture model

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

How to estimate

value of λ ?

unigram $(w_e | s, w)$ $\sum_{s=1}^S = (w_e | w)$

Methodology $(w_e | s, w)$ $\sum_{s=1}^S =$

Train P_1, P_2 on corpus A = "training set"

Fix P_1 and P_2

Estimate λ on corpus C = "development set"

Choose λ to maximize likelihood $\prod_e P_m(w_e | w_{e-1})$ on corpus C.

$$(w_e | w) \lambda + (w_e | w) (1-\lambda) =$$

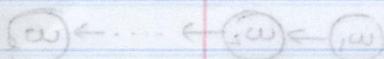
* What happens if we estimate λ on corpus A ?
This would yield $\lambda = 0$ always.

* corpus B = "testing set"

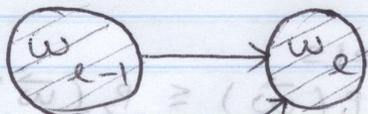
Estimating λ on corpus B is cheating!

QAD* How to estimate λ ? ($\hat{\lambda}$)

Reformulate a mixture model as hidden variable model where $\lambda \in [0, 1]$ appears as value in some hidden node's CPT.



* Belief network



$$P(w_e | w_{e-1}, z) = \begin{cases} P_1(w_e) & \text{if } z=1 \\ P_2(w_e | w_{e-1}) & \text{if } z=2 \end{cases}$$

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

Hidden variable z

Observed variables w_e, w_{e-1}

In this model:

$$P(w_e | w_{e-1}) = \sum_{z=1}^2 P(w_e, z | w_{e-1}) \quad \text{marginalization}$$

$$= \sum_{z=1}^2 P(w_e | z, w_{e-1}) P(z | w_{e-1}) \quad \text{product rule}$$

$$= \sum_{z=1}^2 P(w_e | z, w_{e-1}) P(z) \quad \text{conditional independence}$$

$$= \lambda P_1(w_e) + (1-\lambda) P_2(w_e | w_{e-1})$$

Estimating λ on corpus B is creating a corpus B = "testing set".