



FUDAN SOE SUMMER SCHOOL INTRODUCTION TO AI

Probabilistic Reasoning and Decision Making
Day 8 – Maximum Likelihood Learning

Estimating Parameters of a Bayes Net

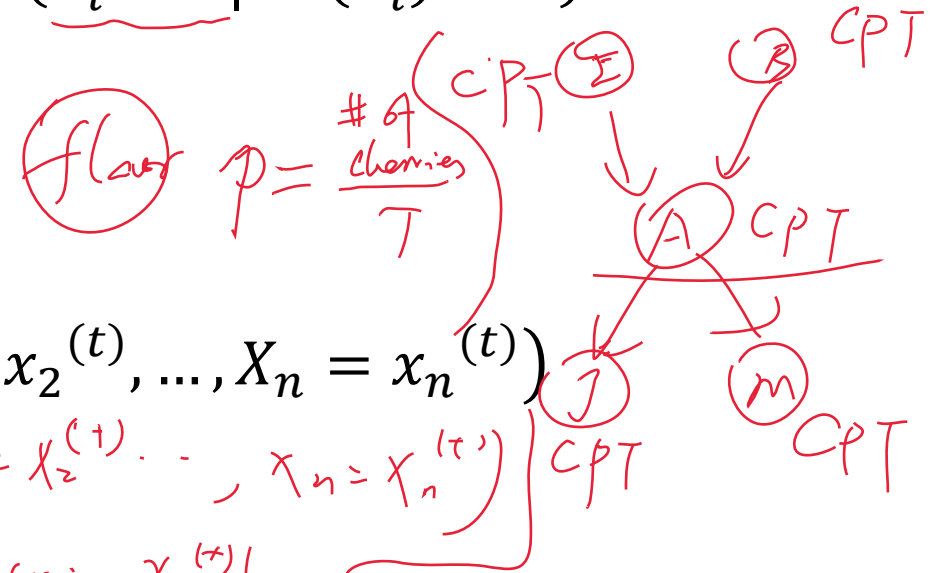
Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

$$\mathcal{L} = \log P(\text{data}) \stackrel{iid}{=} \log \prod_{t=1}^T P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)})$$

$$= \sum_{t=1}^T \log (P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)}))$$

$$P(X_1 = x_1^{(t)}, X_2 = x_2^{(t)}, \dots, X_n = x_n^{(t)}) = \prod_{i=1}^n P(X_i = x_i^{(t)} | Pa(X_i) = \pi_i^{(t)})$$



$$\mathcal{L} = \sum_{t=1}^T \sum_{i=1}^n \log P(X_i = x_i^{(t)} | Pa(X_i) = \pi_i^{(t)})$$

sum over all samples

sum over all variables

E, B	$P(A=1 E, B)$
0 1	—
1 0	—
0 0	—
1 1	—

Estimating Parameters of a Bayes Net

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

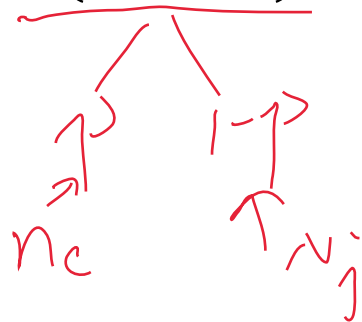
Log-likelihood $\mathcal{L} = \sum_{i=1}^n \sum_{t=1}^T \log P(x_i^{(t)} | pa_i^{(t)})$

Where do we go from here?

Reminder: Single node Bayes Net example:

$$\mathcal{L} = \sum_{t=1}^T \log P(X = x^{(t)})$$

How did we simplify this expression?



$$\frac{\partial \mathcal{L}}{\partial p} = 0$$

$$p = \frac{n_c}{n_c + n_j} = \frac{n_c}{T}$$

Estimating Parameters of a Bayes Net

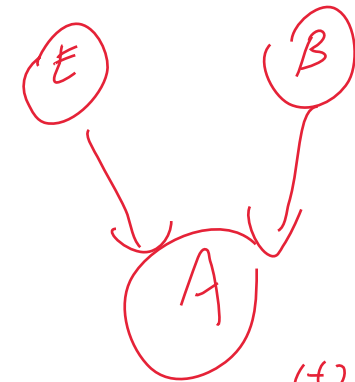
Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

Log-likelihood $\mathcal{L} = \sum_{i=1}^n \sum_{t=1}^T \log P(x_i^{(t)} | pa_i^{(t)})$

$$P(A=0 | E, B=\sim)$$

Where do we go from here?
(example)



data	E	B	A
	0	0	0
	0	0	1
	1	0	0
	1	0	0
	1	1	0
	;	;	;

T rows

define: $\text{Count}_{\pi} (x_i = x_i^{(t)}, Pa(x_i) = \pi_i^{(t)})$

Count(A=0, Pa(A)=00) = 1 Count(A=1, Pa(A)=11) = 0
 Count(A=0, Pa(A)=11) = 1

Let count $(X_i = x_i^{(t)}, Pa(X_i) = Pa_i^{(t)})$ be # of samples where X_i is $x_i^{(t)}$, and parents of X_i is $Pa_i^{(t)}(\pi)$

$$L = \sum_{i=1}^n \left[\sum_{t=1}^T \log P(X_i = x_i^{(t)} | Pa(X_i) = Pa_i^{(t)}) \right]$$

finite # of values

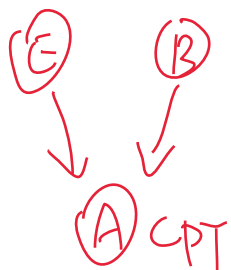
$$= \sum_{i=1}^n \sum_{x_i} \sum_{\pi_i} \text{count}(X_i = x_i, Pa(X_i) = \pi_i) \log P(X_i = x_i | Pa(X_i) = \pi_i)$$

Sum over variables

over the values of X_i

over the value of $Pa(X_i)$

$$T=6, N=3$$



t	E	B	A
1	0	0	0
2	0	0	0
3	0	1	1
4	0	1	0
5	1	1	1
6	1	1	1

$$L = \sum_{i=1}^3 \sum_{t=1}^6 \log P(X_i = x_i^{(t)} | P(x_i) = \tau_i)$$

$$= \sum_{t=1}^6 \log P(E = e^{(t)}) + \sum_{t=1}^6 \log P(B = b^{(t)})$$

$$+ \sum_{t=1}^6 \log P(A = a^{(t)} | \underline{E, B} = \tau_0)$$

$P_A(A)$

$$\log P(A=0 | E=B=00) + \log P(A=0 | E=B=00) + \log P(A=1 | E, B=01)$$

$$+ \log P(A=0 | E, B=01) + \log P(A=1 | E, B=11) + \log P(A=1 | E, B=11)$$

$$2 \log P(A=0 | E, B=00) \dots$$

count(A=1, E, B=11)

$$+ \boxed{2} \log P(A=1 | E, B=11)$$

$$+ 1 \cdot \log P(A=1 | E, B=01) + \log P(A=0 | E, B=01)$$

Estimating Parameters of a Bayes Net

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

Can take the derivative to maximize. Which term are we optimizing over?

- A. $\text{count}(X_i = x, pa_i = \pi)$
- ☒ B. $P(X_i = x | pa_i = \pi)$
- C. $X_i = x$
- D. $pa_i = \pi$
- E. t

Log-likelihood

Explain this in your own words

$$\mathcal{L} = \sum_{i=1}^n \sum_{t=1}^T \log P(x_i^{(t)} | pa_i^{(t)})$$

$$= \sum_{i=1}^n \sum_x \sum_{\pi} \text{count}(X_i = x, pa_i = \pi) \log P(X_i = x | pa_i = \pi)$$

constant

$$\frac{\partial \mathcal{L}}{\partial P(X_i = x | pa_i = \pi)} = 0$$

ML Parameters for Bayes Nets:

if w/o parents
 $P_{ML}(X_i = x) = \frac{\text{count}(X_i = x)}{T}$

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

if w/ parents
 $P(X_i = x | Pa(X_i) = \pi) = \frac{\text{count}(X_i = x, Pa(X_i) = \pi)}{\text{count}(Pa(X_i) = \pi)}$

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$

all possible values for X_i

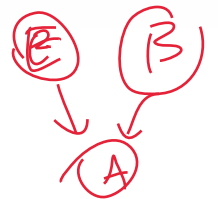
isolate a pattern for parents, find the ratio of a particular value against all possible values for X_i

$$P(A=1 | E B=00) = \frac{0}{2}$$

t	E	B	A
1	0	0	0
2	0	0	0
3	0	1	1
4	0	1	0
5	1	1	1
6	1	1	1

A: 0.5
 B: 0
 C: 1

$$P(A=1 | EB=01) = \frac{1}{1+1} = \frac{1}{2} = .5$$



ML Parameters for Bayes Net: Example

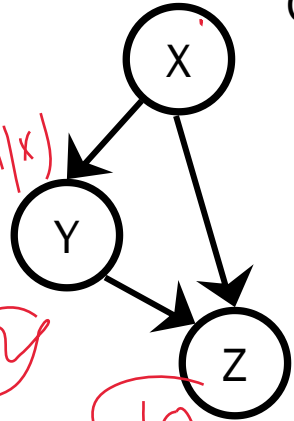
$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$

Observed data:

X	Y	Z
0	0	1
0	1	0
0	1	1
0	1	0
1	0	0
1	0	0
0	1	1
1	0	0
0	1	1
0	0	1
0	1	1
1	0	0

Which of the following is a parameter we would like to estimate?

- ☒ A. $P(X=1)$
- ☐ B. $P(Y=1)$
- ☐ C. $P(X=1|Y=1)$
- ☐ D. More than one of these
- ☐ E. None of these



X, Y and Z are Boolean variables

X	Y	$P(Z=1 X,Y)$
0	0	—
0	1	—
1	0	—
1	1	—

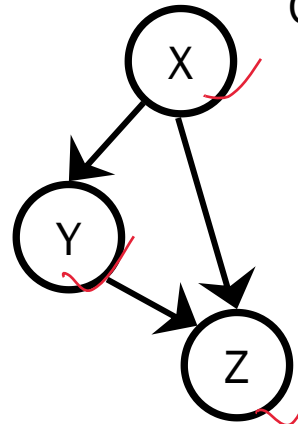
$x/P(Y=1|x)$
0/1

$x/P(x)$
1/1

T: 12 n: 3

4

ML Parameters for Bayes Net: Example



Observed data:

X	Y	Z
0	0	1
0	1	0
0	1	1
0	1	0
1	0	0
1	0	0
0	1	1
1	0	0
0	1	1
0	0	1
0	1	1
1	0	0

X, Y and Z
are Boolean
variables

Not including complements (e.g. $P(X=1)$ and $P(X=0)$), how many different parameters are there to estimate?

- A. 3
- B. 4
- C. 5
- D. 7**
- E. more than 7

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$

Handwritten calculations:

$$P(X=1) \checkmark = \frac{\text{count}(X=1)}{12} = \frac{4}{12} = \frac{1}{3}$$

$$0 = \frac{0}{4} = P(Y=1 | X=1) \quad P(Y=1 | X=0) = \frac{6}{8} = \frac{3}{4}$$

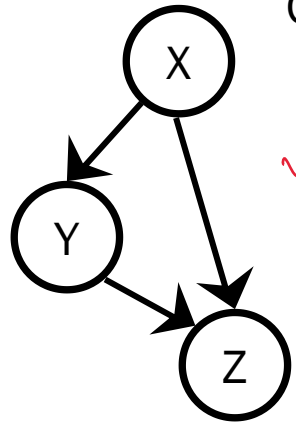
$$1 = P(Z=1 | X=0, Y=0), P(Z=1 | X=0, Y=1)$$

$$P(Z=1 | X=1, Y=0), P(Z=1 | X=1, Y=1)$$

undefined

$$= \frac{0}{\sim} = 0 \quad \frac{4}{6} = \frac{2}{3}$$

ML Parameters for Bayes Net: Example



X, Y and Z
are Boolean
variables

Observed data:

X	Y	Z
0	0	1
0	1	0
0	1	1
0	1	0
1	0	0
1	0	0
0	1	1
1	0	0
0	1	1
0	0	1
0	1	1
1	0	0

What is the ML estimate for $P(Z=1|X=0, Y=0)$?

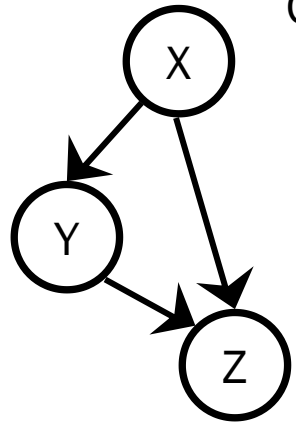
- A. 0
- B. 1/6
- C. 1/2
- D. 1
- E. None of the above

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$

$$= \frac{\text{count}(Z=1, X=0, Y=0)}{\text{count}(X=0, Y=0)} = \frac{2}{2} = 1$$

ML Parameters for Bayes Net: Example

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$



Observed data:

X	Y	Z
0	0	1
0	1	0
0	1	1
0	1	0
1	0	0
1	0	0
0	1	1
1	0	0
0	1	1
0	0	1
0	1	1
1	0	0

X, Y and Z
are Boolean
variables

Which parameter has an undefined ML estimate?

A. $P(X=1)$

B. $P(Y=1 | X=0)$

C. $P(Z=1 | X=0, Y=0)$

D. $P(Z=1 | X=1, Y=1) = \frac{\text{count}(Z=1, X=1, Y=1)}{\text{count}(X=1, Y=1)}$

E. More than one of the above

undefined

0

Summary: Estimating Parameters of a Bayes Net

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

Data set: $\{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$

Log-likelihood

Explain this in your own words

$$\begin{aligned}\mathcal{L} &= \sum_{i=1}^n \sum_{t=1}^T \log P(x_i^{(t)} | pa_i^{(t)}) \\ &= \sum_{i=1}^n \sum_x \sum_{\pi} \text{count}(X_i = x, pa_i = \pi) \log P(X_i = x | pa_i = \pi)\end{aligned}$$

Review: ML Parameters for Bayes Nets

Parameters to estimate (CPTs of the network): $P(X_i = x | Pa(X_i) = \pi)$

$$P_{ML}(X_i = x | pa_i = \pi) = \frac{\text{count}(X_i = x, pa_i = \pi)}{\sum_{x'} \text{count}(X_i = x', pa_i = \pi)}$$

$$= \frac{\text{count}(X_i = x, pa_i = \pi)}{\text{count}(pa_i = \pi)} \quad \text{w/ parents}$$

$$P_{ML}(X_i = x) = \frac{\text{count}(X_i = x)}{\sum} \quad \text{w/o parents}$$

Using the Indicator function I

$$I(x, y) = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{otherwise} \end{cases}$$

- Which of the following correctly expresses the ML estimate

$$\underline{P_{ML}(X_i = x)} = \frac{\text{Count}(X_i = x)}{T} = \frac{1}{T} \left(\sum_{t=1}^T I(x_i^{(t)}, x) \right)$$

$\text{Count}(X_i = x)$

- A. $I(x_i, x) \times T$
- B. $\sum_{i=1}^n I(x_i, x)$
- ☒ C. $\frac{1}{T} \sum_{t=1}^T I(x_i^{(t)}, x)$
- D. $\frac{1}{T} \sum_{t=1}^T I(x_i^{(t)}, x) P(X_i = x_i^{(t)})$
- E. None of these

Suppose we have a belief network with nodes X_1, \dots, X_n , and let $\text{Pa}(X_i)$ denote the parents of node X_i .

A fully-observed dataset for this model can be written as

$$\text{data} = \{(x_1^{(t)}, x_2^{(t)}, \dots, x_n^{(t)})\}_{t=1}^T$$

The log-likelihood of this model is $\mathcal{L} = \log P(\text{data})$. Show that the log-likelihood can be written as

$$\mathcal{L} = \sum_{i=1}^n \sum_x \sum_{\pi} \text{count}(X_i = x, \text{Pa}(X_i) = \pi) \log P(X_i = x | \text{Pa}(X_i) = \pi)$$

Let the graph be as follows:

Nodes: A, B, C, D, E, F

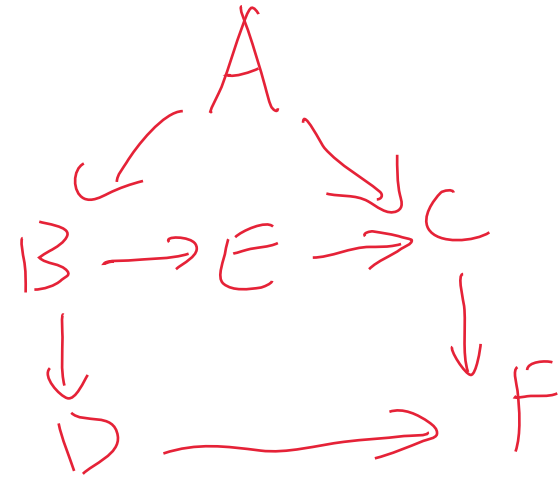
Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

(a) $P(A = a)$

$$= \frac{1}{T} \sum_{t=1}^T \underbrace{I(a_t, a)}$$

↳ for each sample,
is $a_t = a$??



Let the graph be as follows:

Nodes: A, B, C, D, E, F

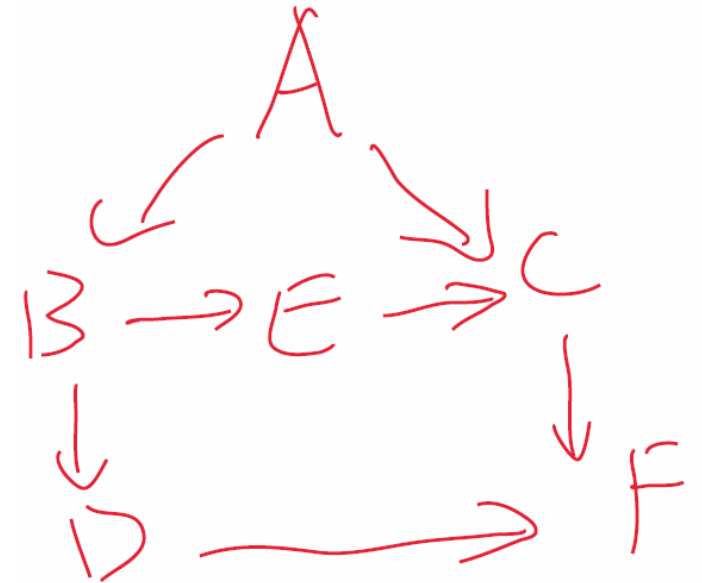
Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

(b) $P(B = b | A = a)$

$$= \frac{\text{count}(A=a, B=b)}{\text{count}(A=a)}$$

$$= \frac{\sum_{t=1}^T I(a_t, a) \cdot I(b_t = b)}{\sum_{t=1}^T I(a_t, a)}$$



Let the graph be as follows:

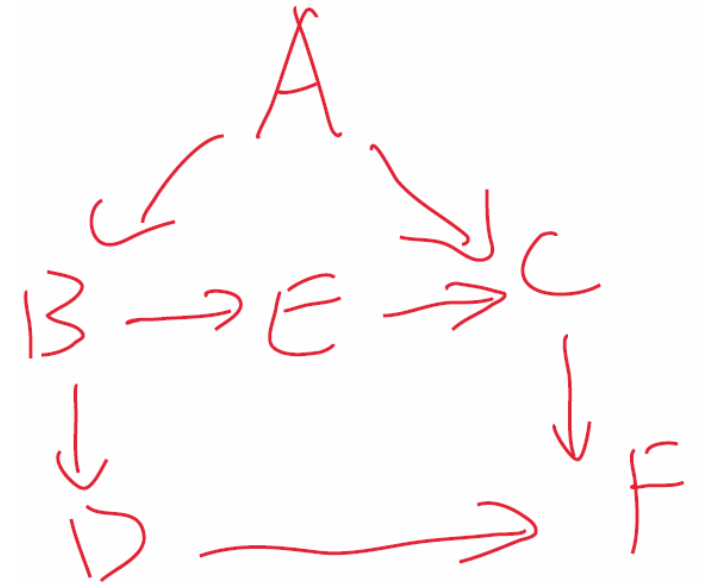
Nodes: A, B, C, D, E, F

Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

(c) $P(C = c | A = a, E = e)$

$$= \frac{\sum_{t=1}^T I(a_t, a) I(e_t, e) I(c_t, c)}{\sum_{t=1}^T I(a_t, a) I(e_t, e)}$$



Let the graph be as follows:

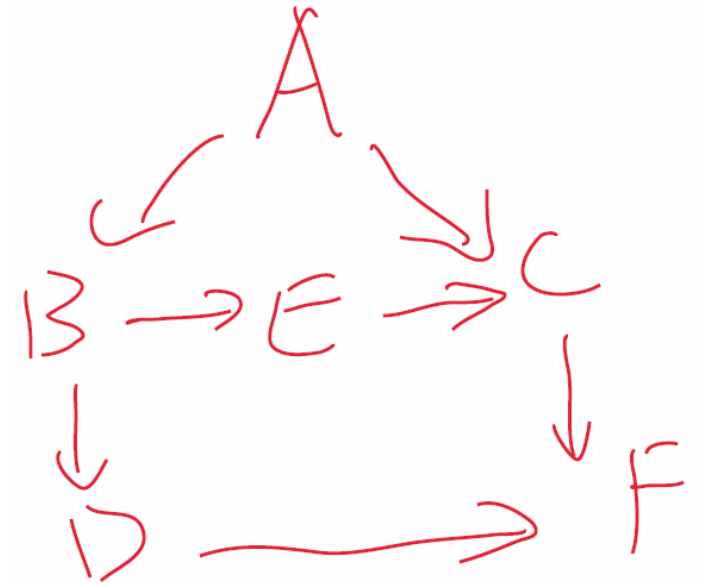
Nodes: A, B, C, D, E, F

Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

(d) $P(D = d | B = b)$

$$= \frac{\sum_t I(d_t, d) I(b_t, b)}{\sum_t I(b_t, b)}$$



Let the graph be as follows:

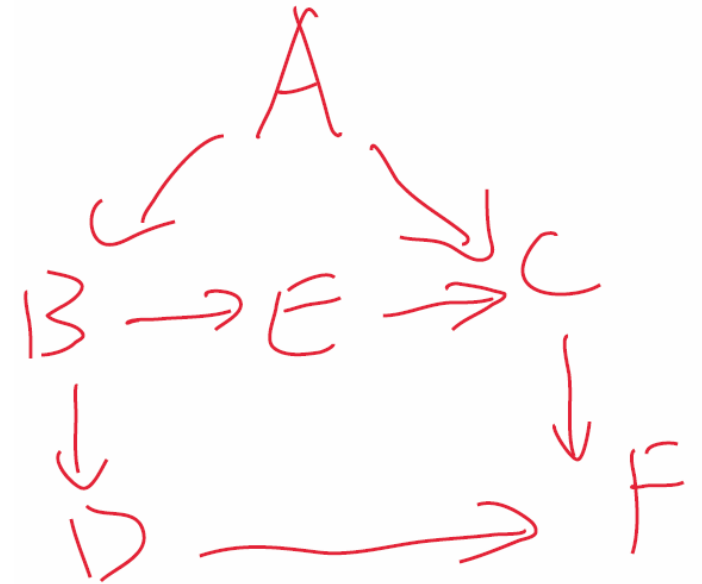
Nodes: A, B, C, D, E, F

Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

(e) $P(E = e | B = b)$

$$= \frac{\sum_t I(b_t, b) I(e_t, e)}{\sum_t I(b_t, b)}$$



Let the graph be as follows:

Nodes: A, B, C, D, E, F

Edges: $A \rightarrow B, A \rightarrow C, B \rightarrow D, B \rightarrow E, E \rightarrow C, C \rightarrow F, D \rightarrow F$

Consider a complete data set of T i.i.d. examples: $\{(a_t, b_t, c_t, d_t, e_t, f_t)\}_{t=1}^T$. Compute the maximum likelihood estimates of the conditional probability tables (CPTs) shown below for this data set. Complete these CPT estimates in terms of the indication function denoted by $I(a, a_t) = 1$ if $a = a_t$ and 0 otherwise.

$$(f) \quad P(F = f | C = c, D = d)$$

$$= \frac{\sum_t I(f_t, f) I(c_t, c) I(d_t, d)}{\sum_t I(c_t, c) I(d_t, d)}$$

Naïve Bayes and Markov Models: Two kinds of Bayes Nets

- A Naïve Bayes model is a Bayes net with a single parent and many children.
- Example: Document Classification

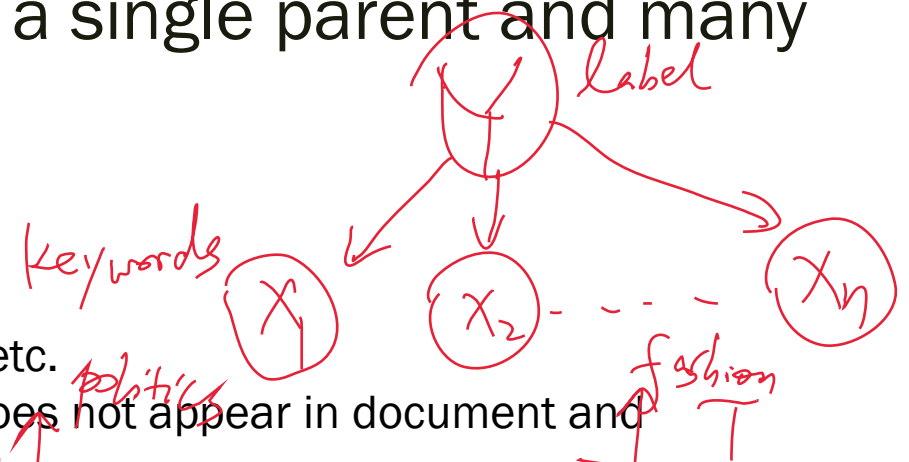
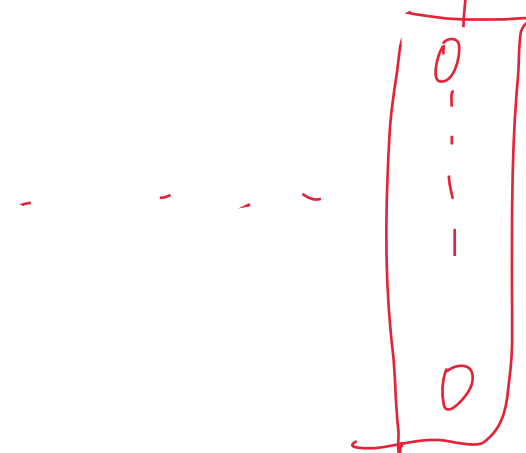
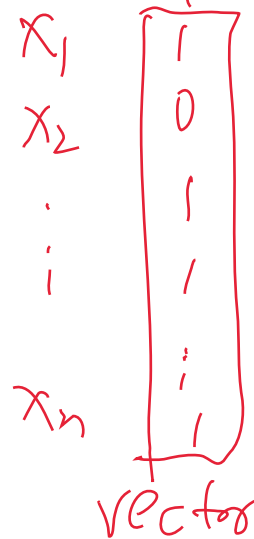
$Y \in \{1, 2, 3, \dots, k\}$ Where 1=sports, 2=fashion, 3=politics, ..., etc.

$X_i \in \{0, 1\}$ for $i=1\dots n$, where 0 means i th word in dictionary does not appear in document and 1 means it does

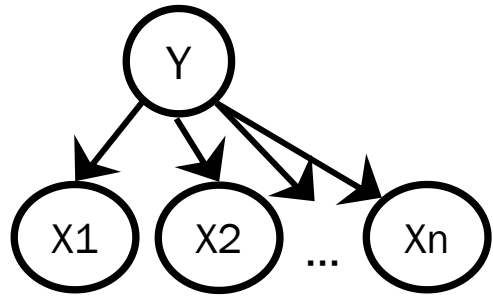
new vector



(label: ?)



Naïve Bayes: Learning and Classification



Learning:

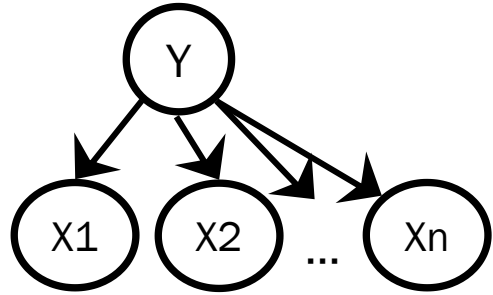
$$P(Y = y) = \frac{\sum_t \mathbb{1}(y_t, y)}{T}$$

$$P(X_i = 1 | Y = y) = \frac{\sum_t \mathbb{1}(y_t, y) \mathbb{1}(x_t, 1)}{\sum_t \mathbb{1}(y_t, y)}$$

Classification

$$\begin{aligned}
 & P(Y = y \mid x_1 = x_1^{(new)}, x_2 = x_2^{(new)}, \dots, x_n = x_n^{(new)}) \\
 &= \frac{P(Y = y, x_1 = x_1^{(new)}, x_2 = x_2^{(new)}, \dots, x_n = x_n^{(new)})}{\sum_{y'} P(Y = y', x_1 = x_1^{(new)}, x_2 = x_2^{(new)}, \dots, x_n = x_n^{(new)})}
 \end{aligned}$$

Naïve Bayes: Learning and Classification



Learning:

$P(Y = y)$: Proportion of documents labeled as category y

$P(X_i = 1|Y = y)$: Proportion of category y documents where word X_i occurs.

Classification

$$\frac{P(Y = y) \prod_{i=1}^n P(X_i = x_i|Y = y)}{\sum_{y'} P(Y = y') \prod_{i=1}^n P(X_i = x_i|Y = y')}$$

y politics : .5
fashion : .2
sports : .1

Strengths and weaknesses?

label : best probability :

Markov Models: Sequential models where each element depends on only elements before it

Example: Language

Let W_i denote the i^{th} word in a sentence. $P(w_1, w_2, w_3, \dots, w_L) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$

Two simplifying assumptions:

1. Finite Context
2. Position Invariance