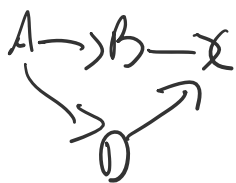


- Last time, What are Bayes Nets?

- Graphical model: V = random variables (features)

E = conditional dependence relationships

Conditional probability table for each variable



$$\Rightarrow P(A, B, C, D) = P(A) P(B|A) P(C|A) P(D|B, C)$$

Classification

- class label (random variable), $c_j \in C$

specific class \uparrow random variable C

- Set of features $X = \langle x_1, x_2, \dots, x_n \rangle$, evidence

- Max a posteriori hypothesis

$$C_{MAP} = \underset{j}{\operatorname{argmax}} P(c_j | x_1, x_2, \dots, x_n)$$

\uparrow posterior distribution

$$= \underset{j}{\operatorname{argmax}} \frac{P(x_1, x_2, \dots, x_n | c_j) P(c_j)}{P(x_1, x_2, \dots, x_n)}$$

for normalization, so we can ignore it

$$= \arg \max_j P(x_1, x_2, \dots, x_n | c_j) P(c_j)$$

↑
Probability of
evidence given
class label

↑
Prior probability
of observing c_j

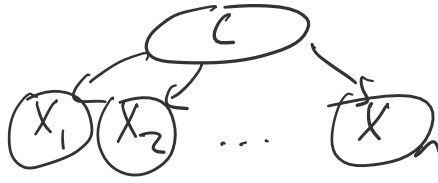
- $P(c_j)$ is not always available
- Assume uniform in that case
- Maximum Likelihood estimate

$$C_{MLE} = P(x_1, x_2, x_3, \dots, x_n | c_j)$$

Typically, we stick with MAP hypothesis

- $C_{MAP} = P(x_1, x_2, \dots, x_n | c_j) P(c_j)$
- How hard is this to calculate in practice?
- $P(c_j)$ is easy!
 - Use frequencies of c_j to estimate $P(c_j)$
- $P(x_1, x_2, \dots, x_n | c_j)$ > Super hard u
 - $O(|X|^n \cdot |C|)$ # of parameters
 - Superexponential growth in terms of data
 - Curse of dimensionality
- What do?

- Naïve Bayes assumption: treat all features $x_i \in \langle x_1, x_2, x_3, \dots, x_n \rangle$ as independent given C



$P(x, y | z) = P(x | z) P(y | z)$ if x, y are conditionally independent given z

Now!

$$C_{MAP} = P(c_j) \prod_i P(x_i | c_j) \leftarrow \text{usually much easier to calculate}$$

- Do suffer information loss
- In practice, works well enough

Application: text classification

- input: some text, label (topic, or genre, or author)
- goal: use text as evidence to predict label
- Step 1, find text features
 - Word frequency
 - n-gram frequency: The quick brown fox
 - pros: temporal relation bigrams:

info
The-quick, quick-brown, brown-fox
 - con: data is sparse

- Position \leftarrow Brent way!
Don't be like Brent

- Each position is a variable

x_1	x_2	x_3	x_4
"	"	"	"
The	Quick	Brown	Fox

- Doesn't scale

- Unreasonable in practice

- Get rid of positional assumptions

- Frequency of words are important!

- Bag of words model

- Text classification depends solely on how often certain words are used